

IntechOpen

Numerical Simulations in Engineering and Science

Edited by Srinivas P. Rao



NUMERICAL SIMULATIONS IN ENGINEERING AND SCIENCE

Edited by **Srinivas P. Rao**

Numerical Simulations in Engineering and Science

<http://dx.doi.org/10.5772/68125>

Edited by Srinivas P. Rao

Contributors

Khurram Mehboob, Mohammad Aljohani, Alexander Khoperskov, Sergey Khrapov, Lei-Yong Jiang, Noreen Sher Akbar, D Tripathi, Z.H Khan, Hiroshi Yokoyama, Akiyoshi Iida, Aleksandr Namgaladze, Maria Knyazeva, Mikhail Karpov, Oleg Zolotov, Roman Yurik, Matthias Foerster, Boris E Prokhorov, Oleg Martynenko, Kenichi Masuda, Dai-Heng Chen, Jean-Luc Autran, Daniela Munteanu, Yury Vasilyevich Yanilkin, Vasily Sofronov, Zhmailo Vadim, Marilou Cadatal-Raduban, Pham Hong Minh, Luong Viet Mui, Nobuhiko Sarukura, Nguyen Dai Hung, J.A. Castro Castro, Bruno Amaral Pereira, Roan Sampaio Souza, Elizabeth Mendes Oliveira,IVALDO LEÃO, Gabino Torres-Vega, Armando Martínez-Pérez, José Trinidad Guillen Bonilla, Héctor Guillén-Bonilla, Antonio Casillas Zamora, Gustavo Adolfo Vega Gómez, Nancy Elizabeth Franco Rodríguez, Alex Guillen Bonilla, Juan Reyes, Despoina Karadimou, Nicolas Markatos, Janja Kramer Stajko, Renata Jecl, Jure Ravnik, Melinda Kovacs, EMOKE DALMA KOVACS, Iordana Astefanoaei, Alexandru Stancu, . Ch V K N S N Moorthy, V Srinivas ., Alexandre De Macêdo Wahrhaftig, Reyolando M. L. R. F. Brasil, Lázaro S. M. S. C. Nascimento, Jorge Oliveira

© The Editor(s) and the Author(s) 2018

The rights of the editor(s) and the author(s) have been asserted in accordance with the Copyright, Designs and Patents Act 1988. All rights to the book as a whole are reserved by INTECHOPEN LIMITED. The book as a whole (compilation) cannot be reproduced, distributed or used for commercial or non-commercial purposes without INTECHOPEN LIMITED's written permission. Enquiries concerning the use of the book should be directed to INTECHOPEN LIMITED rights and permissions department (permissions@intechopen.com). Violations are liable to prosecution under the governing Copyright Law.



Individual chapters of this publication are distributed under the terms of the Creative Commons Attribution 3.0 Unported License which permits commercial use, distribution and reproduction of the individual chapters, provided the original author(s) and source publication are appropriately acknowledged. If so indicated, certain images may not be included under the Creative Commons license. In such cases users will need to obtain permission from the license holder to reproduce the material. More details and guidelines concerning content reuse and adaptation can be found at <http://www.intechopen.com/copyright-policy.html>.

Notice

Statements and opinions expressed in the chapters are those of the individual contributors and not necessarily those of the editors or publisher. No responsibility is accepted for the accuracy of information contained in the published chapters. The publisher assumes no responsibility for any damage or injury to persons or property arising out of the use of any materials, instructions, methods or ideas contained in the book.

First published in London, United Kingdom, 2018 by IntechOpen
eBook (PDF) Published by IntechOpen, 2019

IntechOpen is the global imprint of INTECHOPEN LIMITED, registered in England and Wales, registration number: 11086078, The Shard, 25th floor, 32 London Bridge Street
London, SE19SG – United Kingdom
Printed in Croatia

British Library Cataloguing-in-Publication Data
A catalogue record for this book is available from the British Library

Additional hard and PDF copies can be obtained from orders@intechopen.com

Numerical Simulations in Engineering and Science
Edited by Srinivas P. Rao

p. cm.

Print ISBN 978-1-78923-450-3

Online ISBN 978-1-78923-451-0

eBook (PDF) ISBN 978-1-83881-283-6

We are IntechOpen, the world's leading publisher of Open Access books Built by scientists, for scientists

3,550+

Open access books available

112,000+

International authors and editors

115M+

Downloads

151

Countries delivered to

Our authors are among the
Top 1%

most cited scientists

12.2%

Contributors from top 500 universities



WEB OF SCIENCE™

Selection of our books indexed in the Book Citation Index
in Web of Science™ Core Collection (BKCI)

Interested in publishing with us?
Contact book.department@intechopen.com

Numbers displayed above are based on latest data collected.
For more information visit www.intechopen.com



Meet the editor



Dr. Srinivas P. Rao is presently working as a professor in the Department of Mechanical Engineering at the Institute of Aeronautical Engineering, Hyderabad. In 2008, Dr. Rao earned his PhD degree in Computational Fluid Dynamics. He worked as a scientist in the Division of Computational Fluid Dynamics at the Scientific Engineering and Computing Group (SECG) of the Centre for Development of Advanced Computing (C-DAC), Pune. Dr. Rao has delivered over 15 plenary and keynote lectures or invited lectures at international conferences, taught over 24 courses on CFD and turbulence modeling and computational aerodynamics, and advised 4 doctoral research fellows and research visitor for 36 graduate students (MTech). He has been an editor of international scientific journals and a reviewer for more than 20 journals of Asia, Europe, and the USA. Dr. Rao's research has involved applications of the CFD to the problems of mechanical and aerospace systems, computational physics, turbulence modeling, and recently the application of computational fluid dynamics to biology and medicine.

Contents

Preface XIII

Section 1 Theory, Numerics and Applications 1

Chapter 1 **The Global Numerical Model of the Earth's Upper Atmosphere 3**

Namgaladze Aleksandr, Knyazeva Maria, Karpov Mikhail, Zolotov Oleg, Martynenko Oleg, Yurik Roman, Foerster Matthias and Prokhorov Boris

Chapter 2 **Numerical Simulation of Electronic Systems Based on Circuit-Block Partitioning Strategies 23**

Jorge dos Santos Freitas de Oliveira

Chapter 3 **Numerical Simulation of Fission Product Behavior Inside the Reactor Containment Building Using MATLAB 43**

Khurram Mehboob and Mohammad Subian Aljohani

Chapter 4 **Study of the Numerical Diffusion in Computational Calculations 65**

Despoina P. Karadimou and Nikos-Christos Markatos

Chapter 5 **Simulation of Natural Convection in Porous Media by Boundary Element Method 79**

Janja Kramer Stajnko, Renata Jecl and Jure Ravnik

Chapter 6 **Numerical Simulations of a High-Resolution RANS-FVDM Scheme for the Design of a Gas Turbine Centrifugal Compressor 97**

Chellapilla V K N S N Moorthy and Vadapalli Srinivas

- Section 2 Modeling and Simulation of Problems of Fundamental Physics 115**
- Chapter 7 **Susceptibility of Group-IV and III-V Semiconductor-Based Electronics to Atmospheric Neutrons Explored by Geant4 Numerical Simulations 117**
Daniela Munteanu and Jean-Luc Aufran
- Chapter 8 **Ultrashort Pulse Generation in Ce:LiCAF Ultraviolet Laser 135**
Marilou Cadatal-Raduban, Minh Hong Pham, Luong Viet Mui, Nguyen Dai Hung and Nobuhiko Sarukura
- Chapter 9 **Exact Finite Differences for Quantum Mechanics 163**
Armando Martínez-Pérez and Gabino Torres-Vega
- Chapter 10 **Twin-Grating Fiber Optic Sensors Applied on Wavelength-Division Multiplexing and Its Numerical Resolution 179**
José Trinidad Guillen Bonilla, Héctor Guillen Bonilla, Antonio Casillas Zamora, Gustavo Adolfo Vega Gómez, Nancy Elizabeth Franco Rodríguez, Alex Guillen Bonilla and Juan Reyes Gómez
- Chapter 11 **Numerical Modeling of Chemical Compounds' Fate and Kinetics in Living Organisms: An Inverse Numerical Method for Rate Estimation from Concentration 197**
Kovacs Eموke Dalma and Kovacs Melinda Haydee
- Section 3 Applied Computational Fluid Mechanics 215**
- Chapter 12 **Benchmarks for Non-Ideal Magnetohydrodynamics 217**
Sofronov Vasily, Zhmailo Vadim and Yanilkin Yury
- Chapter 13 **A Numerical Simulation of the Shallow Water Flow on a Complex Topography 237**
Alexander Khoperskov and Sergey Khrapov
- Chapter 14 **Failure Analysis of a High-Pressure Natural Gas Heat Exchanger and its Modified Design 255**
Lei-Yong Jiang, Yinghua Han, Michele Capurro and Mike Benner
- Chapter 15 **Numerical Simulation of Nanoparticles with Variable Viscosity over a Stretching Sheet 273**
Noreen Sher Akbar, Dharmendra Tripathi and Zafar Hayat Khan

- Chapter 16 **Numerical Study of Turbulent Flows and Heat Transfer in Coupled Industrial-Scale Tundish of a Continuous Casting Material in Steel Production 289**
Jose Adilson de Castro, Bruno Amaral Pereira, Roan Sampaio de Souza, Elizabeth Mendes de Oliveira andIVALDO Leão Ferreira
- Chapter 17 **Modeling of the Temperature Field in the Magnetic Hyperthermia 305**
Iordana Astefanoaei and Alexandru Stancu
- Section 4 Vibration Modeling and Numerical Acoustics 325**
- Chapter 18 **Direct and Hybrid Aeroacoustic Simulations Around a Rectangular Cylinder 327**
Hiroshi Yokoyama and Akiyoshi Iida
- Chapter 19 **Analytical and Mathematical Analysis of the Vibration of Structural Systems Considering Geometric Stiffness and Viscoelasticity 349**
Alexandre de M. Wahrhaftig, Reyolando M. L. R. F. Brasil and Lázaro S. M. S. C. Nascimento
- Chapter 20 **Collapse Load for Thin-Walled Rectangular Tubes 371**
Kenichi Masuda and Dai-Heng Chen

Preface

Computational methods have become an indispensable tool in the exploration and analysis of a wide range of physical phenomena. This book offers a detailed exposition of the various numerical methods used in engineering and science; while dealing with various fields and methods, the emphasis is given to the practical and application aspects of such methods. It is an attempt to give the profound background in numerical methods and computer simulation techniques.

This book also illustrates an overview of a broad spectrum of research in computational science and compiles the contents of numerical physics. The chapters focus on the research topics of mathematical modeling and numerical simulation of practical problems of engineering and physics. The diversity of chapters is from challenging physics problems, viz., quantum mechanics and nuclear reactions, to core engineering applications like vibrations, magnetohydrodynamics, and nanotechnology.

This book will be extremely useful for practicing engineers, physicists, and graduate students of science and engineering. The chapters are arranged in such a way that the readers will be able to select the topics appropriate to their interest and need.

I would like to express my gratitude to all those who provided support, offered comments, and assisted in the editing and design. I would also like to thank my family, friends, and students who supported me in this journey of editing the book.

Dr. Srinivas P. Rao
Institute of Aeronautical Engineering
Dundigal, Hyderabad, Telangana, India

Theory, Numerics and Applications

The Global Numerical Model of the Earth's Upper Atmosphere

Namgaladze Aleksandr, Knyazeva Maria,
Karpov Mikhail, Zolotov Oleg, Martynenko Oleg,
Yurik Roman, Foerster Matthias and
Prokhorov Boris

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.71139>

Abstract

The global numerical first principle 3D model of the upper atmosphere (UAM) for the heights 60–100,000 km is presented. The physical continuity, motion, heat balance and electric potential equations for the neutral, ion and electron gases and their numerical solution method are described. The numerical grids, spatial and time integration steps are given together with the boundary and initial conditions and inputs. Testing and obtained geophysical results are given for many observed situations at various levels of solar, geomagnetic and seismic activity.

Keywords: global numerical model, UAM, equations, solution method, grid, integration steps, input, output, upper atmosphere, thermosphere, ionosphere, plasmasphere, magnetosphere, space weather, solar and geomagnetic activity

1. Introduction

The upper atmosphere is a part of the gas envelope of the Earth. It is located from the height $h \approx 60$ km to several R_E of the geocentric distance, where $R_E = 6371$ km is the Earth's radius. It is characterized by a sharp transition from the predominance of the neutral particles to the charged particles.

The upper atmosphere is divided into several height regions depending on its gas composition and dominating physical process: the thermosphere (from ~80–90 km to ~400–800 km) and the exosphere (above ~400–800 km) in relation to the neutral particles; the ionosphere, the plasmasphere and the magnetosphere in relation to the charged particles.

The ionosphere is located at the height range from ~ 50 km (the upper part of the middle atmosphere—mesosphere) to ~ 1000 km. The plasmasphere is located above the ionosphere up to the plasmapause—the geomagnetic force line with $L \sim (2.5 - 7)$, where L is the geocentric distance to the geomagnetic line top expressed in units of R_E . The magnetosphere is the region below the magnetopause, where the solar wind pressure balances the geomagnetic field one. The magnetopause height is $\sim 10R_E$ geocentric distance on the day side and $\sim 100R_E$ on the night one.

The upper atmosphere state is part of Space Weather. It experiences regular annual, seasonal and diurnal variations as well as disturbances caused by the solar, geomagnetic and lithosphere activities, both globally and locally. Along with the meteorological weather, the Space Weather greatly affects human activity. Such disturbances as geomagnetic storms and substorms, auroras lead to disruptions of the radio communication in the HF range up to its black-out, faulty operation of the navigation satellite systems and electronic onboard equipment of the aircrafts. They generate intense geomagnetic-induced currents in the long conducting lines (transmission facilities, telecommunication, cables, railways and oil and gas pipelines), which create failures of the automatic and relay protection systems and, as a consequence, they cause emergency shutdowns of power supply systems, etc. Therefore, monitoring and forecasting of the upper atmosphere state is an extremely important task.

Monitoring of the near-Earth space is conducted by measurements of different physical plasma parameters (temperatures, concentrations and velocities of neutral and charged components, electric and magnetic fields, etc.) at different heights and areas of the globe, both ground-based and on the board of airplanes, rockets and satellites. Despite the growing level of the technical perfectness, it is impossible to provide stable and global monitoring, and “white areas” still remain over the Earth. These areas are those where measurements cannot be conducted or experience various difficulties (especially over the oceans and near the poles). In such cases, the method of the numerical simulation becomes valuable, and the calculations of the desired parameters fill the “white gaps”. The numerical models use the basic fundamental laws of nature to describe quantitatively the near-Earth environment and/or interpret the measurements.

The Upper Atmosphere Model (UAM), described in this chapter, is a global, 3D, self-consistent, numerical model covering $h = 60 - 100,000$ km. It solves the system of equations describing the basic laws of mass, energy, momentum and electric charge conservation and calculates variations of the neutral (O_2 , N_2 , O, H) and charged (electrons and ions O_2^+ , NO^+ , O^+ , H^+) particles' concentrations; temperatures of neutral gas, ions and electrons; the velocities of the particles as well as the electric field potential for any geo- and heliophysical conditions. The UAM takes into account the non-coincidence of the Earth's geographic and geomagnetic axes.

The UAM was developed previously in Kaliningrad under the supervision of Prof. A.A. Namgaladze as the Global Self-consistent Model of the Thermosphere, Ionosphere and Protonosphere (GSM TIP) [1, 2] and further improved in Murmansk. The modern version called as UAM [3, 4] differs from the GSM TIP by implementation of algorithms for the integration with variable latitude steps and by including empirical models of the ionosphere and thermosphere to use their data for initial and boundary conditions and for the UAM testing.

Most of the modern numerical models (NCAR TIE-GCM [5], CTIM [6], CTIPe [7], GITM [8], SWMF [9, 10]) either cover the near-Earth space in very limited ranges in heights and latitudes, or perform a combination of several models without physical coupling between the parts. The UAM covers the near-Earth space as a coupled system and is still superior to all existing models by spatial coverage and resolution. This makes the UAM suitable to investigate a variety of physical processes, both globally and locally.

2. Basic equations

The fundamental conservation laws are used in the UAM: the quasi hydrodynamic equations of the continuity (1), momentum (2), heat balance (3) and electric current continuity (4).

$$\partial n_\alpha / \partial t + \nabla(n_\alpha \vec{v}_\alpha) = Q_\alpha - L_\alpha \quad (1)$$

$$\left(\partial \vec{v}_\alpha / \partial t + [\vec{\Omega} \times [\vec{\Omega} \times \vec{r}]] + 2[\vec{\Omega} \times \vec{v}_\alpha] \right) = \vec{F}_\alpha \quad (2)$$

$$\vec{F}_\alpha = -\nabla p_\alpha + n_\alpha m_\alpha \vec{g} - \sum_i \mu_{ai} v_{ai} n_\alpha (\vec{v}_\alpha - \vec{v}_i) + e n_\alpha (\vec{E} + [\vec{v}_\alpha \times \vec{B}]) \quad (3)$$

$$\rho_\alpha c_{v\alpha} \partial T_\alpha / \partial t + p_\alpha \nabla \vec{v}_\alpha = \nabla(\lambda_\alpha \nabla T_\alpha) + P_{Q\alpha} - P_{L\alpha} + P_{T\alpha}$$

$$\text{div} \vec{j} = 0. \quad (4)$$

where n_α is the concentration of the α -th gas, \vec{v}_α is its velocity vector, Q_α and L_α are its production and loss rates, respectively, ρ_α is the mass density of the α -th gas, $\vec{\Omega}$ is the angular velocity of the Earth's rotation, \vec{r} is the radius-vector directed from the center of the Earth, \vec{F}_α is the force vector acting on the unit gas volume, $p_\alpha = n_\alpha k T_\alpha$ is the partial pressure of the α -th gas, T_α is its temperature, k is Boltzmann's constant, \vec{g} is the vector of the gravity acceleration, m_α is the mass of the α -th gas particle, μ_{ai} is the reduced mass of colliding particles, v_{ai} is the collision frequency, $(\vec{v}_\alpha - \vec{v}_i)$ is the relative velocity of the α -th and i -th gases, e is the elementary charge, \vec{E} and \vec{B} are the electric and magnetic fields, respectively, $c_{v\alpha}$ is the specific heat per unit volume of the α -th gas, λ_α is the thermal conductivity, $P_{Q\alpha}$ is the rate of heat of the α -th gas, $P_{L\alpha}$ is its rate of cooling, $P_{T\alpha}$ is the rate of heat exchange between the α -th gas and other particles and \vec{j} is the electric current density.

Due to the changes in gas composition and the predominance of different physical processes at different heights, the UAM is divided into several computational blocks: the neutral atmosphere and lower ionosphere block; the F2-layer and plasmasphere block and the magnetosphere block. A special block calculates the electric field potential by solving the continuity equation for \vec{j} . Each block covers a particular height range, uses its own particular equations, coordinate system and calculates its own set of parameters and exchanges with other blocks at each time step of numerical integration.

The block of the neutral atmosphere and lower ionosphere uses the spherical geomagnetic coordinate system and calculates the 3D variations of n_n for O, O₂ and N₂ gases, n_e , T_n and \vec{v}_n , where

\vec{v} is the neutral wind velocity vector, within the height range from the model's lower boundary to 520 km, as well as T_i and T_e up to 175 km.

Concentrations $n(\text{N}_2)$ are obtained from the barometric law, $n(\text{O})$ and $n(\text{O}_2)$ from the continuity Eq. (1), which takes view as:

$$\partial n_n / \partial t + \nabla [n_n (\vec{v} + \vec{v}_{dn})] = Q_n - L_n, \quad (5)$$

where \vec{v} is the neutral wind velocity, \vec{v}_{dn} is the diffusion velocity vector, which has only a vertical component equal to the sum of molecular and turbulent diffusion velocities, Q_n and L_n are the production and loss rates, respectively, taking into account the photodissociation of the molecular oxygen by the solar radiation and recombination of O and O₂ in the triple collisions and radiative recombination.

The meridional and zonal components of \vec{v} are obtained from the projection of Eq. (2) on the horizontal axes of the geographical coordinate system:

$$\rho (\partial \vec{v} / \partial t + (\vec{v}, \nabla) \vec{v} + 2[\vec{\Omega} \times \vec{r}])_{hor} = -(\nabla p)_{hor} - \vec{R}_{ni\ hor} + \eta (\nabla^2 \vec{v})_{hor}, \quad (6)$$

where ρ is the mean neutral density, $p = \sum_n n_n k T_n$ is the mean partial pressure, $\vec{R}_{ni\ hor} = \sum_n \sum_i \mu_{ni} v_{ni} n_n (\vec{v} - \vec{v}_i)_{hor}$ is the horizontal projection of the neutral-ion friction force, η is the coefficient of viscosity, n and i lowercases stand for the neutral particles and ions, respectively, v_{ni} is the frequency of collisions between the neutrals and ions, $(\vec{v} - \vec{v}_i)$ is the relative velocity of the neutrals and ions.

Neutral density ρ is calculated from the projection of Eq. (2) on the radii vector \vec{r} directed from the center of the Earth to the particular point in space, taking into account only the gravity force and vertical component of ∇p , that is, fulfilling the hydrostatic equilibrium:

$$\rho g = -\partial p / \partial r, \rho = \sum_n \rho_n = \sum_n n_n m_n, \quad (7)$$

where g is the vertical projection of the sum of the gravitational and centrifugal accelerations.

By summing up for n the continuity Eq. (5) and assuming that the total number of particles remains constant, that is, $\sum_n Q_n = \sum_n L_n$, we obtain:

$$\partial \rho / \partial t + \nabla(\rho \vec{v}) = 0. \quad (8)$$

With the calculated ρ and horizontal components of \vec{v} , we calculate the vertical wind velocity from Eq. (8).

Finally, T_n is calculated from the heat balance Eq. (3) written as:

$$\rho c_v \partial T_n / \partial t + (\vec{v}, \nabla) T_n + p \nabla \vec{v} = \nabla(\lambda_n \nabla T_n) + P_{nQ}^{UV} + P_{nQ}^I + P_{nQ}^C - P_{nL}, \quad (9)$$

where c_v and λ_n are the specific heat per unit volume and thermal conductivity, respectively, P_{nQ}^{UV} , P_{nQ}^I and P_{nQ}^C are the rates of the heating by the ultraviolet (UV) and extra-ultraviolet (EUV) solar radiation, Joule heating due to collisions with ions and heating by particles precipitating from the magnetosphere, respectively; P_{nL} is the rate of the cooling due to radiation.

For concentrations of O_2^+ , N_2^+ and NO^+ ions in the D, E and F1 layers, the photochemical, local heating and heat exchange processes dominate over the transport processes, since the lifetime of the molecular ions is many times smaller than the transport characteristic time due to high collision frequencies of the charged particles with neutrals and each other. Thus, Eq. (1) for the molecular ions in the lower ionosphere block is written as:

$$\partial n_i / \partial t = \partial n(XY^+) / \partial t = Q(XY^+) - L(XY^+), \quad (10)$$

where $n(XY^+) = n_e$ is the total concentration of the molecular ions, $Q(XY^+)$ and $L(XY^+)$ are the production and loss rates due to the ionization by direct UV solar radiation, scatter radiation and ionization by the precipitating electrons, as well as the result of the ion-molecular reactions and dissociative recombination.

The momentum equation for the molecular ions is written as:

$$n_i m_i \vec{g} - \nabla(n_i k T_i) - \sum_n \mu_{in} v_{in} n_i (\vec{v}_i - \vec{v}) + e n_i (\vec{E} + [\vec{v}_i \times \vec{B}]) = 0. \quad (11)$$

In the momentum equation for electrons, we neglect all terms containing m_e due to its smallness:

$$\nabla(n_e k T_e) + e n_e (\vec{E} + [\vec{v}_e \times \vec{B}]) = 0. \quad (12)$$

This gives the electric field of the ambipolar diffusion:

$$\vec{E}_{\parallel} = \nabla(n_e k T_e) / e n_e. \quad (13)$$

The temperatures T_i and T_e are calculated in this block taking into account that up to $h = 175$ km the local heating dominates, and the processes of the heating transport are neglected in the heat balance Eq. (3):

$$3/2 n_i k \partial T_i / \partial t = P_{iQ}^I + P_{iT}^e + P_{iT'}^n \quad (14)$$

$$3/2 n_e k \partial T_e / \partial t = P_{eQ}^P + P_{eQ}^C + P_{eT}^i + P_{eT'}^n \quad (15)$$

where P_{iQ}^I is the rate of Joule's heating, $P_{iT}^e = -P_{iT}^n$ are the rates of the ions' heat exchange with electrons and neutrals, respectively, P_{eT}^i and P_{eT}^n are the rates of the electrons' heat exchange with ions and neutrals, respectively, P_{eQ}^P and P_{eQ}^C are the rates of the electrons' heating by the photoelectrons and electrons precipitating from the magnetosphere, respectively.

In the block of the F2-layer and plasmasphere $n(\text{O}^+)$, $n(\text{H}^+)$, \vec{v}_i , T_i and T_e are calculated for $h = 175 - 100,000$ km. At these heights collision frequencies ν_{in} of charged particles with neutrals are much smaller than ions' gyrofrequencies Ω_i , that is, collisions do not interfere the cyclotron rotation and drift of charged particles. Ions and electrons are fully magnetized, that is, they are tied to the geomagnetic field lines and can move across the lines only under the action of an extraneous force. Because the geomagnetic field has such great influence on the behavior of ions and electrons, we use the magnetic dipole coordinate system in this block. The system of modeling Eqs. (1)–(3), thus, is integrated along the field lines taking into account the electromagnetic plasma drift perpendicular to the magnetic field.

In the continuity equation, along with the production rates Q_i and loss rates L_i of ions O^+ and H^+ by photoionization, corpuscular ionization, chemical reactions between O^+ , O_2 and N_2 , as well as the charge exchange between O^+ and H and between O and H^+ , the transport of ionized components (the divergence of charged particles' flow) is taken into account:

$$D n_i / Dt + \nabla^{par} (n_i v_i^{par}) = Q_i - L_i - n_i \nabla^{per} v^{per}. \quad (16a)$$

Here, superscripts *par* and *per* refer to directions parallel and perpendicular to the geomagnetic field lines, respectively. The operator

$$D/Dt = \partial/\partial t + (v_i^{per}, \nabla), \quad (16b)$$

gives the Lagrangian time derivatives along the trajectory of the charged particle's electromagnetic drift perpendicular to the geomagnetic field line with the velocity:

$$v_i^{per} = v_e^{per} = [\vec{E} \times \vec{B}] / B^2. \quad (16c)$$

The electrons' velocity along the field lines is calculated as:

$$v_e^{par} = \sum_i n_i v_i^{par} / n_e. \quad (17)$$

The atomic ions' motion equation for the i -th atomic ion is given as:

$$\begin{aligned} 2 m_i n_i (\vec{\Omega} \times \vec{v}_i)^{par} &= m_i n_i g^{par} - \nabla^{par} (n_i k T_i) - n_i / n_e \nabla^{par} (n_e k T_e) \\ &\quad - \sum_n \mu_{in} \nu_{in} n_i (v_i^{par} - v^{par}) - \sum_j \mu_{ij} \nu_{ij} n_i (v_j^{par} - v^{par}), \end{aligned} \quad (18)$$

where subscripts i, j, e and n refer to O^+ , H^+ , electrons and neutrals, respectively and g^{par} is the projection of the sum of the gravity and centrifugal accelerations along the geomagnetic field line.

The temperatures T_i and T_e are calculated from the heat balance equations written as:

$$3/2 n_i k (D T_i / Dt + v_i^{par} \nabla^{par} T_i) + n_i k T_i \nabla v_i^{par} - \nabla^{par} (\lambda_i \nabla^{par} T_i) = P_{iQ}^i + P_{iT}^e + P_{iT}^j + P_{iT'}^n \quad (19)$$

$$3/2 n_e k (D T_e / Dt + v_e^{par} \nabla^{par} T_e) + n_e k T_e \nabla v_e^{par} - \nabla^{par} (\lambda_e \nabla^{par} T_e) = P_{eT}^i + P_{eT}^j + P_{eT}^n + P_{eQ}^p + P_{eQ}^c \quad (20)$$

where P_{iQ}^j is the rate of Joule's heating for the ion gas; P_{it}^j , P_{it}^e and P_{it}^n are the rates of the heat exchange between ions, between ions and electrons and between ions and neutrals, respectively; P_{et}^i , P_{et}^j and P_{et}^n are the rates of the electrons' heat exchange with the i -th and j -th ions as well as the neutral gas, respectively; P_{eQ}^p and P_{eQ}^c are the rates of the electrons' heating by the photoelectrons and electrons precipitating from the magnetosphere, respectively.

The magnetospheric block of the UAM calculates the plasma layer n_i from the continuity Eq. (21), \vec{v}_i from the motion Eq. (22), ion pressure p_i from Eq. (23) and the field-aligned current density \vec{j}_{\parallel} (FACs) from Eq. (24) above 175 km:

$$\partial n_i / \partial t + \vec{\nabla} \cdot (n_i \vec{v}_i) = 0, \quad (21)$$

$$e n_i (\vec{E} + [\vec{v}_i \times \vec{B}]) = \vec{\nabla} p_i, \quad (22)$$

$$d(p_i V^\gamma) / dt = 0, \quad (23)$$

$$\vec{j}_{\parallel} = \frac{\vec{e}_z [\vec{\nabla} V \times \vec{\nabla} p_i]}{B}, \quad (24)$$

where V is a half-volume of the geomagnetic field tube, $\gamma = 5/3$ is the adiabatic coefficient for the plasma layer ions and \vec{e}_z is the unit vector along \vec{B} . V is calculated as:

$$V = B \int_0^{\tilde{z}_{max}} B^{-1} dz, \quad (25)$$

where z is a distance along \vec{B} , z_{max} is the distance to the top of the geomagnetic field line.

We assume that the pressure of the plasma layer is isotropic and constant along \vec{B} . The pressure of the plasma layer electrons is negligible in comparison with p_i , that is, the electrons are considered to be cold.

In the block of the electric field potential φ , it is calculated taking into account electric fields of magnetospheric, thermospheric (dynamo), and lower atmospheric (due to the charges transfer into the ionosphere from bellow) origins. The equation of the electric current continuity is written as:

$$\nabla \cdot \vec{j} = \nabla \cdot (\vec{j}_i + \vec{j}_m + \vec{j}_s) = 0, \quad (26a)$$

where \vec{j}_m and \vec{j}_s are the densities of the magnetosphere and lower ionosphere electric currents, and \vec{j}_i is the ionosphere electric current density given by Ohm's law for the plasma:

$$\vec{j}_i = \sigma_{par} \vec{E}_{par} + \sigma_p \vec{E}_{per} + \sigma_H [\vec{B} \times \vec{E}_{per}]. \quad (26b)$$

Here, σ_p and σ_H are the Pedersen and Hall conductivities, respectively, and σ_{par} is the conductivity along \vec{B} .

$$\sigma_p = e^2 n_e \left[\frac{\nu_{in}}{m_i (\Omega_i^2 + \nu_{in}^2)} + \frac{\nu_{en}}{m_e (\Omega_e^2 + \nu_{en}^2)} \right], \quad (26c)$$

$$\sigma_H = e^2 n_e \left[\frac{\Omega_e}{m_e(\Omega_e^2 + \nu_{en}^2)} - \frac{\Omega_i}{m_i(\Omega_i^2 + \nu_{in}^2)} \right], \quad (26d)$$

$$\sigma_{par} = e^2 n_e \left[\frac{1}{m_i \nu_{in}} + \frac{1}{m_e \nu_{en}} \right], \quad (26e)$$

where Ω_i and Ω_e are the ion and electron cyclotron frequencies, respectively.

$$\vec{E}_{par} = \vec{B}(\vec{B}, \vec{E})/B^2, \quad (26f)$$

$$\vec{E}_{per} = \vec{B} \times [\vec{B} \times \vec{E}]/B^2, \quad (26g)$$

where \vec{E} is the electric field vector defined as the sum of the electrostatic field with the potential φ and the induced dynamo field:

$$\vec{E} = -\nabla\varphi + [\vec{v} \times \vec{B}]. \quad (26h)$$

Thus, the equation for the electric current continuity is given as:

$$-\nabla(\hat{\sigma}\vec{E} + \vec{j}_m + \vec{j}_s) = \nabla(\hat{\sigma}(\nabla\varphi - [\vec{v} \times \vec{B}]) - \vec{j}_m - \vec{j}_s) = 0, \quad (27a)$$

where $\hat{\sigma}$ is the tensor of the ionospheric electric conductivity for the coordinate system with axes along \vec{B} , \vec{E} and $[\vec{E} \times \vec{B}]$:

$$\hat{\sigma} = \begin{pmatrix} \sigma_p & 0 & \sigma_H \\ 0 & \sigma_{par} & 0 \\ -\sigma_H & 0 & \sigma_p \end{pmatrix}. \quad (27b)$$

The integration of Eq. (27a) is conducted along the electric conducting layer, from the lower boundary of the UAM up to 175 km, where we neglect the dependency of the electric field intensity on h . Above 175 km, we assume that the plasma is magnetized and the geomagnetic field lines are electrically equipotential.

3. Equations solutions: initial and boundary conditions

The integration of the model equations is carried out using a finite-difference numerical method. The near-Earth environment is considered as a discrete three-dimensional grid, and each computational block uses its own coordinate system. The derivative operators are replaced by the difference ratios, and the solution is obtained in the nodes of the numerical grid. Steps in the numerical grid are chosen depending on the task and the characteristic scale of the investigated process. Usually, the integration step in longitude is chosen as a constant value between 5° and 15°. The integration steps in latitude are variable and depend on the latitude. Near the geomagnetic equator the step is chosen in the range of 2.5–5.0° and 1° near

the poles, because a tighter numerical grid is required for a precise calculation in the aurora zones and across the polar caps. Steps of the numerical integration along the height are also variable. In the blocks where a spherical coordinate system is used, the step is 3 km at the lower boundary of the UAM and increases exponentially with altitude.

The solutions of Eqs. (1)–(4), which are containing derivations of coordinates, require boundary conditions: the distribution of the desired parameters at the boundaries of the numerical blocks. At the lower boundary T_n and n_n are obtained from the empirical model NRLMSISE-00 (further simply MSIS [11]), \vec{v} is calculated from Eqs. (6) and (8) neglecting the viscosity and ion-neutral friction, in the geostrophic approximation, when Eq. (6) contains only the Coriolis acceleration and the pressure gradient of the neutral gas.

For the upper boundary conditions in the block of the neutral atmosphere (at 520 km) the diffusion equilibrium for n_n is used, \vec{v} and T_n are considered to be independent of the height:

$$\partial n_n / \partial r + m_n g / kT = 0, \quad (28)$$

$$\partial v / \partial r = 0, \quad (29)$$

$$\partial T_n / \partial r = 0. \quad (30)$$

In the F2-layer and plasmasphere block the boundary conditions are defined as follows. At the altitude of 520 km, the concentration of the neutral hydrogen is set according to the empirical model of the neutral atmosphere [12]. For atomic ions n_i at the bases of the geomagnetic field lines at 175 km, the production rates are equal to the loss rates:

$$Q_i = L_i. \quad (31)$$

Eqs. (14) and (15) are used for the T_i and T_e at 175 km.

The model equations are integrated along the geomagnetic field lines in the areas with closed lines of force ($\pm 75^\circ$ in geomagnetic latitude). In the regions with open geomagnetic field lines (poleward of 75° geomagnetic latitude) n_p , as well as the heat fluxes of the ions and electrons are set equal to zero at the upper boundary of the UAM. Thus, the model reproduces the condition of the polar wind: a supersonic outflow of plasma from the F2-layer and upper ionosphere of the polar caps along the open lines of the geomagnetic field.

Solutions of those equations, which are containing time derivatives, require initial conditions: parameter distributions at the start of the numerical calculation. For quiet geophysical conditions, stationary solutions are used as initial conditions after the multiple runs of the model. The initial conditions for perturbed periods are the solutions obtained for the previous quiet days. It is also possible to use data from empirical models (NRLMSISE-00 [11], HWM-93 [13], IRI-2001 [14] or IRI-2007 [15]) as initial conditions.

4. Model inputs

The inner state of the simulated space and external forcings acting on it are characterized by the model inputs set up by the user: (1) date and time to set an initial placement of the numerical grid nodes relative to the Sun; (2) spectra of the solar UV and EUV radiation; (3) fluxes of the high energetic electrons precipitating from the magnetosphere; (4) the FACs connecting the ionosphere with the magnetosphere and/or (5) the distribution of φ at the boundaries of the polar cap; (6) the local \vec{j}_s flowing through the lower boundary from below; (7) indices of the geomagnetic activity and (8) components of the interplanetary magnetic field (IMF) and solar wind.

The solar UV and EUV spectra define the coefficients of O_2 dissociation and O_2^+ , N_2^+ , NO^+ and O^+ production rates due to the photoionization of the corresponding neutral components. The UV and EUV spectra dependence on the solar activity is set up according to Ref. [16]. The intensity of the night scatter radiation intensity is 5 kR for the emission with wave length $\lambda = 121.6$ nm and 5 R for the rest emission lines (102.6, 58.4 and 30.4 nm).

The precipitating electrons' fluxes are set up at the upper boundary of the thermosphere, at 520 km, and their intensity I is written as:

$$I(\Phi, \Lambda, E) = I_m(E) \exp\left[-(\Phi - \Phi_m(E))^2/\Delta\Phi(E)^2 - (\Lambda - \Lambda_m(E))^2/\Delta\Lambda(E)^2\right], \quad (32)$$

$$\Phi_m = (\Phi_{md} + \Phi_{mn})/2 + \cos\Lambda(\Phi_{md} - \Phi_{mn})/2, \quad (33)$$

where Φ and Λ are the geomagnetic latitude and longitude, respectively ($\Lambda=0$ corresponds to the midday magnetic meridian); E is the energy of the precipitating electrons; $I_m(E)$ is the maximum intensity of the precipitating electron flux; $\Delta\Phi$ and $\Delta\Lambda$ are the half-widths of the maximum precipitations in latitude and longitude; Φ_{md} and Φ_{mn} are the magnetic latitudes of the maximum precipitations at the midday and midnight meridians, respectively. Specific values for the precipitating parameters in Eqs. (32) and (33) are taken from the empirical models [17, 18].

The magnetospheric sources \vec{j}_m in Eq. (27a) specify the distribution of FACs. The FACs of the Region 1 (R1) flow from the magnetosphere into the ionosphere on the dawn side and out of the ionosphere on the dusk side, at latitudes higher $\pm 75^\circ$. The FACs of the Region 2 (R2) flow in the opposite direction, at areas equatorward of the R1 currents. The distribution of current densities depends on the geophysical conditions and is setup in the UAM in several different ways depending on the task, either as distribution of the FACs in the R1, 2 and the cusp region according to the model [19], or as the distribution of electric potential at the boundary of the polar cap [20] with the FACs in the cusp region and the R2.

The so-called seismogenic electric currents are the vertical electric currents switched on to simulate the ionosphere effects of various mesoscale phenomena in the lower ionosphere, such as earthquakes, thunderstorms, etc. Used as a model input, the vertical \vec{j}_s are added locally to Eq. (27a) at the lower boundary of the UAM.

The geomagnetic activity is used in the UAM by setting up the planetary geomagnetic indices Kp and Ap , indicating global geomagnetic field disturbance, and aurora indices AE , AL and AU , where AU and AL indicate, respectively, the largest increase and lowest decrease of the Northern component of the geomagnetic field in comparison to the background (quiet) value, AE is the sum of AL and AU and characterizes the largest scale of the magnetic field during the substorm in the high-latitude regions.

5. The UAM versions

The UAM provides the possibility of integrating various empirical models and data of the upper atmosphere. The comparison of the self-consistent UAM version and the UAM versions with different combinations of the empirical models allows testing both the UAM and the empirical models.

In the UAM-MSIS version, n_n and T_n are calculated directly from the MSIS [11]. The thermospheric circulation is calculated by the numerical solution of Eqs. (6) and (8) where ∇p from the MSIS is used. Finally, the MSIS is used to set up T_n and n_n at the lower boundary and initial conditions as well.

In the UAM-HWM version the distribution of the horizontal thermospheric wind is calculated using the empirical model HWM-93 [13]. The vertical component of the wind velocity is calculated by the numerical solution of the continuity equation for ρ (Eq. (8)). The HWM-93 is used for the set of the initial conditions and for comparison of the theoretical model winds with observations.

The magnetospheric block of the UAM simulates the transport processes in the plasma sheet by solving the system of the equations for the plasma sheet ions (see item 2). In Ref. [21], the initial values are taken as $p_i = 0,4$ nPa and $n_i = 0,4$ cm⁻³, correspondingly. The program produces more or less realistic p_i distribution and R2 FACs. The problem is that the obtained solution is not stable, and it falls apart after approximately 1 h.

There are several ways to set up FACs spatial-temporal distributions in the UAM, such as empirical data from the magnetic field measurements from the Dynamics Explorer 2 [22] and the Magsat satellites [23], the FACs empirical models by Papitashvili [24] in [25, 26], by Lukianova [27] in Ref. [28] and MFACE [29] in Ref. [30]. All these versions with various FACs take into account the dependence of FACs on the interplanetary magnetic field (IMF). Such methods of setting the FACs distribution allow using any other empirical data of FACs.

In the UAM version [31], the positions of the auroral oval boundaries, the values of electron flux intensities and the latitudes and longitudes of the intensity maxima were set from precipitation patterns observed by DMSP. The spectra of the precipitating electrons are assumed to be Maxwellian in this case.

The UAM-P version [32, 33], created in Potsdam, differs from the UAM by the electric field block simulation. This block uses magnetic dipole coordinates instead of spherical geographical

ones within $h = 80 - 526$ km. It is assumed that \vec{E} does not change along \vec{b} inside the ionospheric current layer. After the integration, it is also represented as a 2D-distribution. It allows to exclude the conductivity parallel to the magnetic field and to keep the vertical electric field inside the current-carrying layer. As a result, the lower latitudinal and equatorial electric field distributions as well as the current system of these areas look more correctly.

The Canadian Ionosphere and Atmosphere Model (Canadian IAM or C-IAM) is comprised of the extended Canadian Middle Atmosphere (CMAM) and the UAM, currently coupled in a one-way manner [34]. This version was used to investigate the physical mechanisms responsible for forming the four-peak longitudinal structure of the 135.6 nm ionospheric emission observed by the IMAGE satellite over the tropics at 20:00 local time from March 20 to April 20, 2002. The study showed that main mechanism is driven by the diurnal eastward propagating tide with zonal number 3.

6. Results

During its development and improvement, the UAM was used to perform a number of numerical experiments aimed at testing and comparison of calculation results with measurements and other models. The UAM successfully showed its ability to reproduce the general behavior of ionospheric and thermospheric parameters such as at low and high geomagnetic and solar activity conditions. A good agreement of the numerical simulations' results was achieved in comparison with observations by incoherent scatter radars located at various latitudes and longitudes (ISRs) [31], digital ionosonde CADI in the Voeykovo Main Geophysical Observatory [35], several chains of the ionospheric tomographical receivers [36–38], satellites, including CHAMP and GPS [3, 4, 26, 30, 39–42], as well as with empirical models such as various types of IRI, MSIS, HWM [42–47], etc. and other theoretical models [42].

Numerical modeling of the upper atmosphere behavior during substorms. These relatively short disturbances of the geomagnetic field have duration from about 0.5 to 3 h. They are generated at the magnetopause and in the magnetosphere tail via the magnetic field lines reconnection processes and connected with the polar upper atmosphere in the auroral zones via FACs including the current wedges. The substorm auroral currents are reflected by the auroral magnetic activity indexes AL , AU and AE . The UAM takes into account these indexes. This allows the modeling of the upper atmospheric behavior during substorms via the UAM simulations. The results were presented in [4, 48–53] including the cusp and auroral oval behavior, energetic magnetospheric electron precipitations, electric fields, current wedge and internal atmospheric gravity waves generation. The main role of the thermospheric heating due to the soft electron precipitation was shown for the thermosphere substorm effects.

Numerical modeling of the upper atmosphere behavior during geomagnetic storms. These geomagnetic disturbances are global and have duration of about several days. Usually they include several substorms on the background of the geomagnetic field depression created by the ring current development. The geomagnetic activity indexes Dst and Kp characterize the geomagnetic storms. These indexes are included in the UAM inputs as well. This allows to model the upper atmospheric behavior during magnetic storms via the UAM simulations.

The calculation results were presented in [25, 26, 30, 36–38, 54–59] including the main ionospheric trough dynamics due to the ionosphere-magnetosphere convection and non-coincidence of the geographical and geomagnetic axes of the Earth. The physical mechanisms of the negative and positive F2-layer ionospheric storms (electron concentration decreases and increases, correspondingly) formation were described. The main role of the thermospheric composition (atomic to molecular neutral gas concentration ratio) and winds disturbances in the magnetic storm ionospheric effects was demonstrated in these UAM calculations.

In addition to the UAM testing and comparison with the observations for the different levels of solar and geomagnetic activities, the model has been widely used to study **the ionosphere response on the local sources in the lower atmosphere**, such as disturbances associated with the processes of the earthquakes' preparation processes [26, 60–63]. Numerical UAM calculations showed that the electric fields of 5–10 mV/m effects on the F2-layer plasma by the electromagnetic drift in the crossed geomagnetic field and the electric field of the seismic origin. The vertical electric current, flowing through the lower boundary of the ionosphere with the density of ~ 20 nA/m², is required for the generation of the seismogenic electric field of ~ 10 mV/m [60–64]. The important role of the aerosols over the tectonic faults was underlined in this process due to the very low recombination rate for the charged aerosols.

Thus, the UAM was tested and used in many helio- and geophysical situations. Nevertheless, the amount of the UAM simulations remains to be insufficient despite the IT progress. This is related with the specifics of the geophysics as science at all. The near-Earth space environment varies due to the solar, seismic and human activities. This does not allow performing the repeated fixed experiments as in usual physical laboratories to obtain correctly the standard statistical error estimates. Moreover, the observations themselves are very limited. None of them has 3D spatial and time resolutions satisfying to the requirements of the modern technical means of the Space weather practical usage. This was well known long ago [65] and such models as UAM are aimed solving this important problem.

7. Conclusions

Further development of the UAM means a huge amount of further numerical experiments to its mathematical and physical quality. These experiments have to take into account all modern achievements of the numerical mathematics and computer science. The numerical grids, steps, various iterations, etc. should be tested to find their optimal combinations for the best stability and accuracy. The user's manuals should be constructed, including the UAM website. The UAM prognostic features have to be improved by modeling many case studies for various helio- and geophysical situations especially for geomagnetic and seismogenic disturbances. Comparisons with ground-based and satellite observations, empirical and other theoretical models have to be made continuously. The frame approach should be widely used by including separate observational, empirical and theoretical blocks into the UAM, such as the real geomagnetic field, polar wind, plasma sheet, electric fields, lower atmosphere, aerosols, tides, etc. An international cooperation is absolutely necessary for these future scientific UAM perspectives.

Acknowledgements

We thank many people who have worked with the UAM, developed and improved it, such as A. N. Namgaladze, M. A. Volkov, E. N. Doronina, Yu. V. Romanovskaya, I. V. Korableva, Yu. A. Shapovalova, M. G. Botova, M. I. Rybakov, I. V. Artamonov, E. V. Vasilieva, V. A. Medvedeva, V. A. Shlykov, L. A. Chernyuk and many others.

Author details

Namgaladze Aleksandr^{1*}, Knyazeva Maria¹, Karpov Mikhail^{1,2}, Zolotov Oleg¹, Martynenko Oleg³, Yurik Roman⁴, Foerster Matthias⁵ and Prokhorov Boris⁵

*Address all correspondence to: namgaladze@yandex.ru

1 Murmansk Arctic State University, Murmansk, Russia

2 Immanuel Kant Baltic Federal University, Kaliningrad, Russia

3 York University, Toronto, Canada

4 Polar Geophysical Institute, Murmansk, Russia

5 Helmholtz Centre Potsdam, GFZ German Research Centre for Geosciences, Potsdam, Germany

References

- [1] Namgaladze AA, Korenkov YN, Klimenko VV, Karpov IV, Bessarab FS, Surotkin VA, et al. Global model of the thermosphere-ionosphere-protonosphere system. *Pure and Applied Geophysics*. 1988;**127**(2/3):219-254. DOI: 10.1007/BF00879812
- [2] Namgaladze AA, Korenkov YN, Klimenko VV, Karpov IV, Surotkin VA, Naumova NM. Numerical modeling of the thermosphere-ionosphere-protonosphere system. *Journal of Atmospheric and Solar–Terrestrial Physics*. 1991;**53**(11/12):1113-1124. DOI: 10.1016/0021-9169(91)90060-K
- [3] Namgaladze AA, Martynenko OV, Namgaladze AN. Global model of the upper atmosphere with variable latitudinal integration step. *International Journal of Geomagnetism and Aeronomy*. 1998;**1**(1):53-58
- [4] Namgaladze AA, Martynenko OV, Volkov MA, Namgaladze AN, Yurik RY. High-latitude version of the global numerical model of the Earth's upper atmosphere. *Proceedings of the MSTU*. 1998;**1**(2):23-84

- [5] Qian L, Burns AG, Emery BA, Foster B, Lu G, Maute A, et al. The NCAR TIE-GCM. In: Huba J, Schunk R, Khazanov G, editors. *Modeling the Ionosphere-Thermosphere System*. 1st ed. Chichester, UK: John Wiley & Sons, Ltd; 2014. DOI: 10.1002/9781118704417.ch7
- [6] Fuller-Rowell TJ, Rees D, Quegan S, Moffett RJ, Codrescu MV, Millward GH. A Coupled thermosphere-ionosphere model (CTIM). In: Schunk RW, editor. *STEP Handbook of Ionospheric Models*. 1st ed. Utah, USA: Utah State University and SCOSTEP (Scientific Committee on Solar Terrestrial Physics); 1996. p. 217-238
- [7] Codrescu MV, Negrea C, Fedrizzi M, Fuller-Rowell TJ, Dobin A, Jakowsky N, et al. A real-time run of the Coupled Thermosphere Ionosphere Plasmasphere Electrodynamics (CTIPE) model. *Space Weather*. 2012;**10**(2). DOI: 10.1029/2011SW000736
- [8] Ridley AJ, Deng Y, Toth G. The global ionosphere-thermosphere model. *Journal of Atmospheric and Solar-Terrestrial Physics*. 2006;**68**(8):839-864. DOI: 10.1016/j.jastp.2006.01.008
- [9] Tóth G, Holst B, Sokolov IV, De Zeeuw DL, Gombosi TI, Fang F, et al. Adaptive numerical algorithms in space weather modeling. *Journal of Computational Physics*. 2012;**231**(3):870-903. DOI: 10.1016/j.jcp.2011.02.006
- [10] Tóth G, Sokolov IV, Gombosi TI, Chesney DR, Clauer CR, De Zeeuw DL, et al. Space Weather Modeling Framework: A new tool for the space science community. *Journal of Geophysical Research*. 2005;**110**(A12). DOI: 10.1029/2005JA011126
- [11] Picone JM, Hedin AE, Drob DP, Aikin AC. NRLMSISE-00 empirical model of the atmosphere: Statistical comparisons and scientific issues. *Journal of Geophysical Research*. 2002;**107**(A12). DOI: 10.1029/2002JA009430
- [12] Jacchia LG. Thermospheric temperature, density and composition: New models. *Journal of Geophysical Research*. 1977;**375**:1-106
- [13] Hedin AE, Fleming EL, Manson AH, Schmidlin FJ, Avery SK, Clark RR, et al. Empirical wind model for the upper, middle and lower atmosphere. *Journal of Atmospheric and Terrestrial Physics*. 1996;**58**:1421-1447. DOI: 10.1016/0021-9169(95)00122-0
- [14] Bilitza D. International reference ionosphere 2000. *Radio Science*. 2001;**36**:261-275
- [15] Bilitza D, Reinisch BW. International reference ionosphere 2007: Improvements and new parameters. *Advances in Space Research*. 2008;**42**:599-609
- [16] Ivanov-Kholodny GS, Nusinov AA. Shortwave solar radiation and its impact on the atmosphere and ionosphere (in Russian). *Investigations of Outer Space*. 1987;**26**:80-154
- [17] Fuller-Rowell TJ, Evans DS. Height-integrated Pedersen and hall conductivity patterns inferred from the Tiros-NOAA satellite data. *Journal of Geophysical Research*. 1987;**92**(7):7606-7618

- [18] Hardy DA, Gussenhoven MS, Holeman E. A statistical model of auroral electron precipitation. *Journal of Geophysical Research*. 1985;**90**:4229-4248
- [19] Iijima T, Potemra TA. The amplitude distribution of field-aligned currents at northern high latitudes observed by Triad. *Journal of Geophysical Research*. 1976;**81**:2165-2174
- [20] Weimer DR, Maynard NC, Burke WJ, Liebrecht C. Polar cap potentials and the auroral electrojet indices. *Planetary and Space Science*. 1990;**38**(9):1207-1222. DOI: 10.1016/0032-0633(90)90028-O
- [21] Artamonov IV, Vasilyeva EV, Medvedeva VA, Namgaladze AA. A new version of the magnetospheric block for the global numerical upper atmosphere model. In: *Physics of Auroral Phenomena, XXVI Annual Apatity Seminar, Abstracts*. Apatity, Russia: PGI of RAS; 2003. p. 21
- [22] Artamonov IV, Vasilyeva EV, Namgaladze AA, Martynenko OV. On the electric potential pattern corresponding region 2 field-aligned currents derived from the DE2 satellite measurements. In: *Physics of Auroral Phenomena. Proceedings of the XXIX Annual Seminar*; Apatity, Russia. 2006. p. 71-74
- [23] Artamonov IV, Vasilyeva EV, Martynenko OV, Namgaladze AA. Mathematical simulation of the electric potential distributions using region 2 field-aligned currents obtained from different satellite data. In: *Proceedings of the 6th International Conference "Problems of Geocosmos"*, Saint-Petersburg, Russia. 2006. p. 20-23
- [24] Papitashvili VO, Christiansen F, Neubert T. A new model of field-aligned currents derived from high-precision satellite magnetic field data. *Geophysical Research Letters*. 2002;**29**(14). DOI: 10.1029/2001GL014207
- [25] Förster M, Prokhorov BE, Namgaladze AA, Holschneider M. Numerical modeling of solar wind influences on the dynamics of the high-latitude upper atmosphere. *Advances in Radio Science*. 2012;**10**:299-312. DOI: 10.5194/ars-10-299-2012
- [26] Namgaladze AA, Förster M, Prokhorov BE, Zolotov OV. Electromagnetic drivers in the upper atmosphere: Observations and modeling. In: Bychkov V, Golubkov G, Nikitin A, editors. *The Atmosphere and Ionosphere, Physics of Earth and Space Environments*. 1st ed. Houten, Netherlands: Springer; 2013. p. 165-219. DOI: 10.1007/978-94-007-2914-8_4
- [27] Lukianova R, Christiansen F. Modeling of the global distribution of ionospheric electric fields based on realistic maps of field-aligned currents. *Journal of Geophysical Research*. 2006;**111**:A03213. DOI: 10.1029/2005JA011465
- [28] Knyazeva MA, Namgaladze AA, Beloushko KE. Field-aligned currents influence on the ionospheric electric fields: Modification of the upper atmosphere model. *Russian Journal of Physical Chemistry B*. 2015;**9**(5):758-763. DOI: 10.1134/S1990793115050206
- [29] He M, Vogt J, Lühr H, Sorbalo E, Blagau A, Le G, et al. A high-resolution model of field-aligned currents through empirical orthogonal functions analysis (MFACE). *Geophysical Research Letters*. 2012;**39**:L18105. DOI: 10.1029/2012GL053168

- [30] Prokhorov BE, Förster M, Namgaladze AA, Holschneider M. Using MFACE as input in the UAM to specify the MIT dynamics. *Journal of Geophysical Research, Space Physics*. 2014;**119**(8):6704-6714. DOI: 10.1002/2014JA019981
- [31] Namgaladze AA, Zubova YV, Namgaladze AN, Martynenko OV, Doronina EN, Goncharenko LP, et al. Modelling of the ionosphere/thermosphere behaviour during the April 2002 magnetic storms: A comparison of the UAM results with the ISR and NRLMSISE-00 data. *Advances in Space Research*. 2006;**37**(2):380-391. DOI: 10.1016/j.asr.2005.04.013
- [32] Prokhorov BE, Förster M, Stolle C, Lesur V, Namgaladze AA, Holschneider M. The ionospheric current system and its contribution to the Earth's magnetic field. *Geophysical Research Abstracts*. 2016;**18**:ST4.1, EGU2016-6902
- [33] Prokhorov BE, Förster M, Lesur V, Lesur V, Namgaladze AA, Holschneider M, Stolle C. Using the UAM-P model to specify the additional magnetic field generated by the system of the ionospheric currents. *Geophysical Research Abstracts*. 2017;**19**:EMRP2.1/ESSI1.15/GI1.12, EGU2017-10822
- [34] Martynenko OV, Fomichev VI, Semeniuk K, Beagley SR, Ward WE, McConnell JC, et al. Physical mechanisms responsible for forming the 4-peak longitudinal structure of the 135.6 nm ionospheric emission: First results from the Canadian IAM. *Journal of Atmospheric and Solar - Terrestrial Physics*. 2014;**120**:51-61. DOI: 10.1016/j.jastp.2014.08.014
- [35] Rybakov MV, Namgaladze AA, Karpov MI. Comparison of ionospheric parameters calculated with UAM and measured at Voeykovo observatory. *Geomagnetism and Aeronomy*. 2016;**56**(5):604-609. DOI: 10.1134/S0016793216040186
- [36] Namgaladze AN, Evstafiev OV, Khudukon BZ, Namgaladze AA. Model interpretation of the ionospheric F-region electron density structures observed by ground-based satellite tomography at sub-auroral and auroral latitudes in Russia in January-May 1999. *Annales Geophysicae*. 2003;**21**(4):1005-1016
- [37] Shapovalova YA, Namgaladze AA, Namgaladze AN, Khudukon BZ. Stratification of the main ionospheric trough as an effect of the noncoincidence of the Earth's geomagnetic and geographic axes. In: *Proceedings of the XXVI Annual Apatity Seminar "Physics of Auroral Phenomena"*. PGI; 2003;**26**:87-90
- [38] Korableva IV, Namgaladze AA, Namgaladze AN. High-latitude ionosphere during magnetic storms of October 26, 2003–November 1, 2003: Tomographic reconstructions and numerical modeling. *Geomagnetism and Aeronomy*. 2008;**48**(5):642-651
- [39] Namgaladze AA, Doronina EN, Förster M. The role of electric fields and magnetospheric electron precipitations for the formation of the equatorial total mass density minimum. In: *"Physics of Auroral Phenomena", Proc. 29 Annual Seminar*. 2006. pp. 238-240
- [40] Doronina EN, Namgaladze AA, Förster M. A model interpretation of the CHAMP neutral mass density measurements. In: *Proceedings of the 6th International Conference "Problems of Geocosmos"* St. Petersburg, Russia. 2006. p. 58-61

- [41] Förster M, Namgaladze AA, Doronina EN, Prokhorov BE. High-latitude thermospheric winds: Satellite data and model calculations. *Russian Journal of Physical Chemistry B*. 2011;**5**(3):439-446. DOI: 10.1134/S1990793111030043
- [42] Shim JS, Rastaetter L, Kuznetsova M, Bilitza D, Codrescu M, Coster AJ, et al. CEDAR-GEM Challenge for Systematic Assessment of Ionosphere/Thermosphere Models in Predicting TEC during the 2006 December Storm Event. *Space Weather*. 05 October 2017;**15**:19. DOI: 10.1002/2017SW001649
- [43] Martynenko OV, Namgaladze AA, Namgaladze AN, Shlykov VA. Numerical modeling of the longitudinal variations in the near-earth plasma, In: *Physics of Auroral Phenomena, Proceedings of the XXIII Annual Apatity Seminar, Apatity PGI-00-01-108*. 2001. pp. 49-52
- [44] Knyazeva MA, Namgaladze AA. A model study of the seasonal and solar activity variations of the enhanced electron density regions in the night-time ionospheric F2-layer and plasmasphere. In: *Physics of Auroral Phenomena, Proc. XXIX Annual Seminar, Apatity, Kola Science Centre, Russian Academy of Science*. 2006. pp. 225-228
- [45] Knyazeva MA, Namgaladze AA. An influence of the thermospheric wind variations on the enhanced electron density regions in the night-time ionospheric F2-layer and in the plasmasphere. In: *Proc. 6th International Conference "Problems of Geocosmos", Saint-Petersburg State University*. 2006. pp. 91-94
- [46] Knyazeva MA, Namgaladze AA, Zubova Yu V. Numerical modeling of the Weddel Sea anomaly. In: *Proceedings of the 8th International Conference Problem of the Geocosmos, St. Petersburg Petrodvorets*. 2010. pp. 121-125
- [47] Botova MG, Romanovskaya YV, Namgaladze AA. Latitudinal variations and altitude profiles of Ionospheric parameters: Comparison of theoretical and empirical model results. *Russian Journal of Physical Chemistry B*. 2015;**9**(5):764-769. DOI: 10.1134/S1990793115050164
- [48] Namgaladze AA, Martynenko OV, Namgaladze AN, Volkov MA, Korenkov YN, Klimenko VV, et al. Numerical simulation of an ionospheric disturbance over EISCAT using a global ionospheric model. *Journal of Atmospheric and Terrestrial Physics*. 1996;**58**(1):297-306. DOI: 10.1016/0021-9169(95)00037-2
- [49] Volkov MA, Namgaladze AA. Models of field-aligned currents needful to simulate the substorm variations of the electric field and other parameters observed by EISCAT. *Annales Geophysicae*. 1996;**14**(12):1356-1361. DOI: 10.1007/s00585-996-1356-0
- [50] Leontyev SV, Namgaladze AA, Namgaladze AN, Bogdanov NN. Thermospheric meridional winds in the vicinity of the auroral zone: Observations and modeling. *Journal of Atmospheric and Solar - Terrestrial Physics*. 1998;**60**(2):215-226. DOI: 10.1016/S1364-6826(97)00055-2

- [51] Namgaladze AA, Namgaladze AN, Volkov MA. Numerical modelling of the thermospheric and ionospheric effects of magnetospheric processes in the cusp region. *Annales Geophysicae*. 1996;**14**(12):1343-1355. DOI: 10.1007/s00585-996-1343-5
- [52] Namgaladze AA, Namgaladze AN. Numerical modelling of the thermospheric and ionospheric effects of the soft electron precipitation at the cusp. In: *Proceedings of the 1st Meeting Workshop: Magnetic Reconnection at the Magnetopause and Aurora Dynamics*; Verlag der Osterreichischen Akademie der Wissenschaften. 1996. p. 177-183
- [53] Namgaladze AA, Namgaladze AN, Volkov MA. Seasonal effects in the ionosphere-thermosphere response to the precipitation and field-aligned current variations in the cusp region. *Annales Geophysicae*. 1998;**16**(10):1283-1298. DOI: 10.1007/s00585-998-1283-3
- [54] Namgaladze AA, Yurik RYu. Plasmasphere state effect on the positive phase of the ionospheric storm. In: *Proceedings of the XXVI Annual Apatity Seminar*; Apatity, Russia. 2003. pp. 79-82
- [55] Namgaladze AA, Namgaladze AN, Yurik RY. Global modeling of the quiet and disturbed upper atmosphere. *Physics and Chemistry of the Earth*. 2000;**25**(5-6):533-536. DOI: 10.1016/S1464-1917(00)00071-4
- [56] Namgaladze AA, Foerster M, Yurik RY. Analysis of the positive ionospheric response to a moderate geomagnetic storm using a global numerical model. *Annales Geophysicae*. 2000;**18**(4):461-477. DOI: 10.1007/s00585-000-0461-8
- [57] Namgaladze AA, Namgaladze AN, Fadeeva Yu V, Goncharenko LP, Salah JE. Lower thermosphere and ionosphere behaviour during a strong magnetic storm of March 31, 2001: Modelling and comparison with the Millstone Hill incoherent scatter radar measurements. *Proceedings of the MSTU*. 2003;**6**(1):87-92
- [58] Namgaladze AA, Namgaladze AN, Martynenko OV, Doronina EN, Knyazeva MA, Zubova Yu V. Numerical modeling of the thermosphere, ionosphere and plasmasphere behaviour during the April 2002 magnetic storms. In: *Proceedings of the XXVI Annual Apatity Seminar Preprint PGI*. 2003. pp. 74-78
- [59] Förster M, Namgaladze AA, Yurik RY. Thermospheric composition changes deduced from geomagnetic storm modeling. *Geophysical Research Letters*. 1999;**26**(16):2625-2628. DOI: 10.1029/1999GL900514
- [60] Namgaladze AA, Klimenko MV, Klimenko VV, Zakharenkova IE. Physical mechanism and mathematical modeling of earthquake Ionospheric precursors registered in Total electron content. *Geomagnetism and Aeronomy*. 2009;**49**(2):252-262. DOI: 10.1134/S0016793209020169
- [61] Zolotov OV, Namgaladze AA, Zakharenkova IE, Martynenko OV, Shagimuratov II. Physical interpretation and mathematical simulation of ionospheric precursors of earthquakes at midlatitudes. *Geomagnetism and Aeronomy*. 2012;**52**(3):390-397. DOI: 10.1134/S0016793212030152

- [62] Karpov MI, Namgaladze AA, Zolotov OV. Modeling of Total electron content disturbances caused by electric currents between the earth and the ionosphere. *Russian Journal of Physical Chemistry B*. 2013;7(5):595-598. DOI: 10.1134/S1990793113050187
- [63] Namgaladze AA, Karpov MI. Conductivity and external electric currents in the global electric circuit. *Russian Journal of Physical Chemistry B*. 2015;9(4):754-757. DOI: 10.1134/S1990793115050231
- [64] Namgaladze AA. Earthquakes and global electrical circuit. *Russian Journal of Physical Chemistry B*. 2013;7(5):589-593. DOI: 10.1134/S1990793113050229
- [65] Nisbet JS. On the Construction and Use of the Pennsylvania State MKI Model. *Ionosphere Research Laboratory*; 1970 98 p

Numerical Simulation of Electronic Systems Based on Circuit-Block Partitioning Strategies

Jorge dos Santos Freitas de Oliveira

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.75490>

Abstract

Numerical simulation of complex and heterogeneous electronic systems can be a very challenging issue. Circuits composed of a combination of analog, mixed-signal and digital blocks or even radio frequency (RF) blocks, integrated in the same substrate, are very difficult to simulate as a whole at the circuit level. The main reason is because they contain a lot of state variables presenting very distinct properties and evolving in very widely separated time scales. Examples of practical interest are systems-on-a-chip (SoCs), very common in mobile electronics applications, as well as in many other embedded electronic systems. This chapter is intended to briefly review some advanced circuit-level numerical simulation techniques based on circuit-block partitioning schemes, which were especially designed to address the simulation challenges brought by this kind of circuits into the computer-aided-design (CAD) field.

Keywords: numerical simulation, electronic systems, multirate schemes, circuit-block partition

1. Introduction

Electronic circuit simulation has emerged in the 1970s, triggered by the necessity of engineers having a tool to help in design and analysis of integrated circuits (ICs). Since probing internal nodes of semiconductor chips is extremely difficult, or even prohibitive in almost all cases, manufacturing integrated circuits having no help of a simulation tool would lead to an unbearable set of successive physical prototypes until a final solution was achieved. A simulation package combining device modeling and numerical simulation will help engineers to verify correctness and debug circuits during their design, avoiding physical prototyping and reducing product development expenses. Over the years, the continuous scaling of semiconductor devices

and the increasing complexity of electronic architectures have been making CAD tools more and more important for circuits and systems designing. Demands for continuously providing new systems' functionalities, lower-power consumptions or higher transmission rates (e.g., RF and microwave communication systems) are typical requisites that have led to a complex scenario of highly heterogeneous electronic systems. Conventional algorithms are not capable of simulating such kind of electronic systems in an efficient way. The justification for such ineffectiveness relies on the fact that standard simulation techniques do not perform any distinction between nodes or blocks within the circuits, treating all the variables in the same manner. This causes all the blocks of the circuit (analog, mixed-signal, digital or RF blocks) to be computed with the same numerical scheme, without taking their nature into consideration. To cope with this scenario, some advanced numerical simulation algorithms based on circuit-block partition have been proposed in recent years in the scientific literature. The most important ones are briefly reviewed in this chapter.

2. Review of basic circuit simulation concepts

2.1. Mathematical modeling of electronic systems

Dynamic behavior of electronic systems is modeled as systems of differential algebraic equations (DAEs) involving voltages, currents, charges and fluxes. These systems are usually obtained via nodal analysis, or modified nodal analysis (MNA), which consists of applying the Kirchhoff's current law (KCL) to each electrical node and writing the branch currents in terms of the circuit node voltages using the corresponding constitutive relations to each circuit element. Such systems have, in general, the following form:

$$\mathbf{p}(\mathbf{y}(t)) + \frac{d\mathbf{q}(\mathbf{y}(t))}{dt} = \mathbf{x}(t), \quad (1)$$

in which $\mathbf{x}(t) \in \mathbb{R}^n$ is the vector of independent stimuli (voltage or current sources) to the circuit, $\mathbf{y}(t) \in \mathbb{R}^n$ is the vector of unknowns (voltages and currents waveforms) and n is the total number of unknowns. $\mathbf{p} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ describes the memoryless elements in the circuit (linear and nonlinear elements, as resistors, nonlinear voltage-controlled current sources, etc.), while $\mathbf{q} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ models all linear and nonlinear reactive circuit elements, as capacitors or inductors, represented as voltage-dependent electric charges or current-dependent magnetic fluxes, respectively.

Applying the chain differentiation rule to the reactive term of the DAE system of (1), we are able to obtain

$$\frac{d\mathbf{q}(\mathbf{y})}{d\mathbf{y}} \frac{d\mathbf{y}(t)}{dt} = \mathbf{x}(t) - \mathbf{p}(\mathbf{y}(t)), \quad (2)$$

or,

$$\mathbf{M}[\mathbf{y}(t)] \frac{d\mathbf{y}(t)}{dt} = \mathbf{x}(t) - \mathbf{p}(\mathbf{y}(t)), \quad (3)$$

where $\mathbf{M}[\mathbf{y}(t)]$ is habitually denoted as the *mass matrix*. If $\mathbf{M}[\mathbf{y}(t)]$ is invertible then it is possible to convert (3) into the following *ordinary differential equations'* (ODE) system,

$$\frac{d\mathbf{y}(t)}{dt} = \mathbf{M}[\mathbf{y}(t)]^{-1}(\mathbf{x}(t) - \mathbf{p}(\mathbf{y}(t))), \quad (4)$$

which can be rewritten in the classical form

$$\frac{d\mathbf{y}(t)}{dt} = \mathbf{f}(t, \mathbf{y}(t)), \quad (5)$$

ordinarily utilized in the mathematical literature. When $\mathbf{M}[\mathbf{y}(t)]$ is singular, the DAE system of (3) will not degenerate into a ODE system, but it is often possible to express it as a set of algebraic equations combined with a set of differential equations of the form of (5).

2.2. Classic SPICE-like simulation

The most natural way of simulating the dynamic behavior of an electronic circuit is to numerically time-step integrate, in time domain, the DAE system, or the ODE system, modeling its operation. This means that the solution of (1), or (5), has to be computed over a specified time interval $[t_0, t_{Final}]$ from a specific initial condition $\mathbf{y}(t_0) = \mathbf{y}_0$, leading to the so-called *initial value problems*

$$\mathbf{p}(\mathbf{y}(t)) + \frac{d\mathbf{q}(\mathbf{y}(t))}{dt} = \mathbf{x}(t), \quad \mathbf{y}(t_0) = \mathbf{y}_0, \quad t_0 \leq t \leq t_{Final}, \quad \mathbf{y}(t) \in \mathbb{R}^n, \quad (6)$$

or

$$\frac{d\mathbf{y}(t)}{dt} = \mathbf{f}(t, \mathbf{y}(t)), \quad \mathbf{y}(t_0) = \mathbf{y}_0, \quad t_0 \leq t \leq t_{Final}, \quad \mathbf{y}(t) \in \mathbb{R}^n. \quad (7)$$

Computing the solution of (6), or (7), is frequently referred to as *transient analysis* and can be done by using initial value solvers, such as linear multistep methods (LMM) [1] or one-step methods (the popular Runge-Kutta (RK) schemes) [2, 3]. Either LMM or RK methods offer a large variety of explicit and implicit schemes, with very distinct properties in terms of order (accuracy) and numerical stability.

This time-step integration technique was used in the first digital computer programs of circuit analysis initially developed at the Electronics Research Laboratory of the University of California, Berkeley, in the early 1970s, and is still today the most widely used technique for such purpose. It is the core of all SPICE (which means Simulation Program with Integrated Circuit Emphasis) or SPICE-like computer programs, available in many commercial simulators.

2.3. Periodic steady-state simulation

As described earlier, SPICE-like simulation tools are primarily focused on transient analysis. However, in some cases electronics designers are essentially interested in obtaining circuits' steady-state responses and not their transient regimes. The reason for that is because some specific properties of the circuits are better characterized, or simply only defined, in steady-state (e.g., impedance, voltage gain, current gain, harmonic distortion, signal to noise ratio, etc.).

SPICE tools are not adequate for computing steady-state responses of circuits presenting very different time constants, or high Q resonances, as is typically the case of RF and microwave circuits. This is so because they have to pass through the lengthy process of integrating all transients, and expecting them to vanish. Indeed, in such cases time-step integration can be extremely inefficient, since the number of discretization time steps used by the numerical integration scheme will be dramatically enormous. This is because the time interval over which the differential equations must be numerically integrated is set by the lowest frequency, or by how long the circuit takes to achieve steady-state, while the length of the time steps is constrained by the highest frequency component.

Periodic steady-state response of an electronic circuit is a regime where its unknowns are a set of generic waveforms (node voltages and branch currents) presenting a common period. Computing the periodic steady-state response, without having to first integrate all the transients, consists of finding the initial condition, $\mathbf{y}(t_0)$, for the DAE, or ODE, system describing the circuit's operation, such that the values of the unknowns at the end of one period T match their values at the beginning of that period, that is to say, $\mathbf{y}(t_0) = \mathbf{y}(t_0 + T)$. These problems (evaluating the solution to a differential system that satisfies constraints at two or more distinct time instants) are referred to as *boundary value problems*. In this case, we have a *periodic boundary value problem* that can be formulated as:

$$\mathbf{p}(\mathbf{y}(t)) + \frac{d\mathbf{q}(\mathbf{y}(t))}{dt} = \mathbf{x}(t), \quad \mathbf{y}(t_0) = \mathbf{y}(t_0 + T), \quad t_0 \leq t \leq t_0 + T, \quad \mathbf{y}(t) \in \mathbb{R}^n, \quad (8)$$

or, in the ODE form, as

$$\frac{d\mathbf{y}(t)}{dt} = \mathbf{f}(t, \mathbf{y}(t)), \quad \mathbf{y}(t_0) = \mathbf{y}(t_0 + T), \quad t_0 \leq t \leq t_0 + T, \quad \mathbf{y}(t) \in \mathbb{R}^n, \quad (9)$$

in which $\mathbf{y}(t_0) = \mathbf{y}(t_0 + T)$ is known as the *periodic boundary condition*.

The most widely used techniques for computing the periodic steady-state solution of electronic circuits are briefly reviewed in the following: the shooting-Newton method [4] and the harmonic balance method [5, 6].

2.3.1. The shooting-Newton method

The shooting-Newton method [4] is the time-domain technique most commonly used by the electronic design automation (EDA) community for computing periodic steady-state solutions of electronic circuits. Shooting-Newton is an iterative technique that: (1) starts by computing the solution of the circuit for one period T using a LMM or RK integrator, considering some

guessed initial condition (which is in general determined from a previous DC analysis); (2) then the computed solution at the end of the period is checked, and if it does not agree with the initial condition, the initial condition is wisely updated and (3) the circuit is then re-simulated for one period with the adjusted initial condition, and this process is repeated until the solution at the end of the period matches the initial condition.

Shooting is an iterative solver that uses an initial value technique to solve a boundary value problem. With the purpose of providing some technical details on the implementation of the shooting-Newton method, let us consider (8). Since with shooting, we perform numerical time-step integration of the DAE system from $t = t_0$ until $t = t_0 + T$, the main difficulty is that we do not know a priori for which initial condition $\mathbf{y}(t_0)$ has to be considered that will lead to the steady-state solution, that is, that will satisfy the periodic boundary condition $\mathbf{y}(t_0) = \mathbf{y}(t_0 + T)$. Thus, the key aspect of shooting-Newton relies on finding the solution of

$$\mathbf{y}(t_0) = \mathbf{y}(t_0 + T) \Leftrightarrow \mathbf{y}(t_0) - \mathbf{y}(t_0 + T) = 0. \quad (10)$$

Let us now define $\phi(\mathbf{y}(t_0), T) = \mathbf{y}(t_0 + T)$, where ϕ is known as the *state-transition function* [4, 7], and rewrite (10) as

$$\phi(\mathbf{y}(t_0), T) - \mathbf{y}(t_0) = 0. \quad (11)$$

Although electronic circuits may operate in strongly nonlinear regimes, their state-transition functions are often moderately nonlinear (or even quite linear). This means that slight perturbations on the initial condition (starting state) produce almost proportional perturbations in the subsequent time states. Taking this aspect into account, it is straightforward to conclude that (11) can be iteratively solved in an efficient way with the Newton's method, which will lead to the following iterative scheme:

$$\phi(\mathbf{y}^{[r]}(t_0), T) - \mathbf{y}^{[r]}(t_0) + \left[\frac{\partial \phi(\mathbf{y}^{[r]}(t_0), T)}{\partial \mathbf{y}(t_0)} - I \right] \Bigg|_{\mathbf{y}(t_0) = \mathbf{y}^{[r]}(t_0)} [\mathbf{y}^{[r+1]}(t_0) - \mathbf{y}^{[r]}(t_0)] = 0, \quad (12)$$

where I is the $n \times n$ identity matrix. The cost of the solution of the linear system of (12) is dominated by the computational effort required to evaluate the derivative of the state-transition function (usually referred to as the *sensitivity matrix*). This matrix is computed taking into account the chain differentiation rule, that is, taking into consideration that $\phi(\mathbf{y}(t_0), T)$ is, in fact, the numerical vector \mathbf{y}_K , with K being the total number of time steps in the interval $[t_0, t_0 + T]$, which depends on the previous step value \mathbf{y}_{K-1} , which, itself, depends on \mathbf{y}_{K-2} , and so forth. The sensitivity matrix is then given by

$$\frac{\partial \phi(\mathbf{y}(t_0), T)}{\partial \mathbf{y}(t_0)} = \frac{\partial \mathbf{y}_K}{\partial \mathbf{y}_{K-1}} \cdot \frac{\partial \mathbf{y}_{K-1}}{\partial \mathbf{y}_{K-2}} \cdot \dots \cdot \frac{\partial \mathbf{y}_1}{\partial \mathbf{y}_0}. \quad (13)$$

Although solving (12) and computing the sensitivity matrix (13) involve substantial computational cost, shooting-Newton converges to the steady-state solution much faster than classic time-step integration. This is the reason why it is the time-domain steady-state engine most widely used in the circuit simulation field.

2.3.2. The harmonic balance method

The harmonic balance (HB) method [4–6] is a frequency domain steady-state simulation technique widely used by the EDA community. In contrast to time domain tools, which represent waveforms as a set of time samples, frequency domain techniques represent periodic signals using coefficients in a sum of complex exponentials (or sines and cosines) harmonically related. The main advantage of HB over time-domain techniques (e.g., shooting-Newton) is that it can represent steady-state solutions (voltage and current waveforms) very accurately using a small number of coefficients. This is especially evident for moderately nonlinear circuits excited by smooth waveforms, in which significant reductions in the computational cost are achieved when HB is used as the simulation tool. However, it must be noted that HB is not suitable for dealing with strongly nonlinear regimes producing waveforms with sharp transitions. In such case, the large number of terms required in the Fourier series expansions will make HB very inefficient.

With the purpose of providing a brief and intuitive explanation of the HB method, let us consider (8). For achieving simplicity in our formulation, let us consider a very small circuit driven by a single periodic source $x(t) \in \mathbb{R}$ satisfying $x(t) = x(t + T)$ and whose dynamic behavior is described by a unique unknown, $y(t) \in \mathbb{R}$ (the generalization to $x(t) \in \mathbb{R}^n$, $y(t) \in \mathbb{R}^n$ is straightforward). Given that the steady-state regime of the circuit will be periodic with the same period T , both the stimulus and the steady-state solution can be represented as the Fourier series

$$x(t) = \sum_{k=-\infty}^{+\infty} X_k e^{jk\omega_0 t}, \quad y(t) = \sum_{k=-\infty}^{+\infty} Y_k e^{jk\omega_0 t}, \quad (14)$$

in which $\omega_0 = 2\pi/T$ is the fundamental frequency. If we substitute (14) into (8), and consider an appropriate harmonic truncation at some order $k = K$, we will obtain

$$p\left(\sum_{k=-K}^{+K} Y_k e^{jk\omega_0 t}\right) + \frac{d}{dt} \left[q\left(\sum_{k=-K}^{+K} Y_k e^{jk\omega_0 t}\right) \right] = \sum_{k=-K}^{+K} X_k e^{jk\omega_0 t}. \quad (15)$$

The HB method converts the differential problem of (8) into the frequency domain, obtaining the $(2K + 1)$ algebraic equations system

$$\mathbf{F}(\mathbf{Y}) = \mathbf{P}(\mathbf{Y}) + j\Omega\mathbf{Q}(\mathbf{Y}) - \mathbf{X} = 0 \quad (16)$$

where $\mathbf{Y} = [Y_{-K}, \dots, Y_0, \dots, Y_K]^T$, $\mathbf{X} = [X_{-K}, \dots, X_0, \dots, X_K]^T$, $j\Omega = \text{diag}(-jK\omega_0, \dots, 0, \dots, jK\omega_0)$.

and \mathbf{P} , \mathbf{Q} are vectors containing the Fourier coefficients of $p(y(t))$ and $q(y(t))$, respectively. The algebraic equations system of (16) is usually denoted as the *harmonic balance system*, which can be iteratively solved using the Newton method

$$\mathbf{F}(\mathbf{Y}^{[r]}) + \frac{d\mathbf{F}(\mathbf{Y})}{d\mathbf{Y}} \Big|_{\mathbf{Y}=\mathbf{Y}^{[r]}} [\mathbf{Y}^{[r+1]} - \mathbf{Y}^{[r]}] = 0, \quad (17)$$

in which

$$\frac{d\mathbf{F}(\mathbf{Y})}{d\mathbf{Y}} = \mathbf{J}(\mathbf{Y}) = \mathbf{G}(\mathbf{Y}) + j\Omega\mathbf{C}(\mathbf{Y}) \quad (18)$$

is the $(2K + 1) \times (2K + 1)$ composite conversion matrix, known as the *Jacobian matrix* of the error function $\mathbf{F}(\mathbf{Y})$. \mathbf{G} and \mathbf{C} denote the $(2K + 1) \times (2K + 1)$ conversion matrices (Toeplitz) [7] corresponding to $g(y(t)) = dp(y(t))/dy$ and $c(y(t)) = dq(y(t))/dy$.

3. Univariate time-domain partitioned simulation engines

3.1. Time-domain latency

As highlighted in the Introduction of this chapter, dynamic regimes of operation of some electronic systems may involve signals (voltages and currents) presenting widely distinct time evolution rates. Typical examples of that are coupled analog-digital systems or combined technologies of baseband analog, digital and RF blocks, in the same circuit. In these examples, very fast signals and slowly varying signals cohabit in the same framework. This property, the one of having in the same problem signals presenting rapid time rates of change, while others evolve in a very slow way (or remain approximately constant within a certain time window) is usually denoted as *time-domain latency*. It must be pointed out that time-domain latency is a phenomenon that is not restricted to heterogeneous circuits. For instance, regimes of operation of pure digital electronic systems typically present a set of variables remaining practically constant within a specific time interval, while others evidence quick variations (fast transitions) in that interval.

In Section 2, we have seen that SPICE-like simulation engines (which are based on time-step integration schemes) are widely used for computing the numerical solution of electronic systems. However, when dealing with circuits presenting time-domain latency, that is, containing node voltages and brunch currents evolving at very different rates, traditional SPICE simulators become inefficient because they expend unnecessary work on the computation of the slowly varying components. This is so because classic initial value solvers (as LMM or RK schemes) integrate all the DAE, or ODE, unknowns with the same step size.

3.2. Partitioned algorithms for time-step integration

To deal with the earlier described time-domain latency in an effective way, some modern partitioned algorithms for univariate time-step integration have been proposed in the recent years in the scientific literature [8–10]. These powerful techniques, denoted as multirate Runge-Kutta (MRK) schemes, split the ODE system of (5) into coupled fast and slow (latent) subsystems, obtaining

$$\begin{aligned} \frac{d\mathbf{y}_F(t)}{dt} &= \mathbf{f}_F(t, \mathbf{y}_F, \mathbf{y}_L), & \mathbf{y}_F(t_0) &= \mathbf{y}_{F,0} \\ \frac{d\mathbf{y}_L(t)}{dt} &= \mathbf{f}_L(t, \mathbf{y}_F, \mathbf{y}_L), & \mathbf{y}_L(t_0) &= \mathbf{y}_{L,0} \end{aligned} \quad (19)$$

with

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_F \\ \mathbf{y}_L \end{bmatrix}, \quad \mathbf{y}_F \in \mathbb{R}^{n_F}, \quad \mathbf{y}_L \in \mathbb{R}^{n_L}, \quad n_F + n_L = n, \quad (20)$$

where \mathbf{y}_F is the vector containing the fast-varying signals and \mathbf{y}_L is the vector holding the slowly varying (latent) ones. The former will be integrated with a small step length h (*microstep*), whereas the latter will be integrated with a large step size H (*macrostep*). The number of microsteps within each macrostep is an integer that we will denote by m . Hence, $H = m \cdot h$. The generic formulation of a MRK scheme is given in the following [8, 9].

Consider two Runge-Kutta methods of s and \bar{s} stages represented by their Butcher tableaus [2] $(\mathbf{b}, \mathbf{A}, \mathbf{c})$ and $(\bar{\mathbf{b}}, \bar{\mathbf{A}}, \bar{\mathbf{c}})$,

$$\begin{array}{c|ccc} c_1 & a_{11} & a_{12} & \cdots & a_{1s} \\ c_2 & a_{21} & a_{22} & \cdots & a_{2s} \\ \vdots & \vdots & \vdots & & \vdots \\ c_s & a_{s1} & a_{s2} & \cdots & a_{ss} \\ \hline & b_1 & b_2 & \cdots & b_s \end{array} \quad \begin{array}{c|ccc} \bar{c}_1 & \bar{a}_{11} & \bar{a}_{12} & \cdots & \bar{a}_{1\bar{s}} \\ \bar{c}_2 & \bar{a}_{21} & \bar{a}_{22} & \cdots & \bar{a}_{2\bar{s}} \\ \vdots & \vdots & \vdots & & \vdots \\ \bar{c}_{\bar{s}} & \bar{a}_{\bar{s}1} & \bar{a}_{\bar{s}2} & \cdots & \bar{a}_{\bar{s}\bar{s}} \\ \hline & \bar{b}_1 & \bar{b}_2 & \cdots & \bar{b}_{\bar{s}} \end{array} \quad (21)$$

The MRK method conceived for efficiently computing the numerical solution of the partitioned differential system of (19), using a microstep h for integrating the fast unknowns and a macrostep H for integrating the latent unknowns, is defined as follows [8, 9].

The fast-varying vector is obtained by

$$\mathbf{y}_F(t_0 + \lambda h + h) \simeq \mathbf{y}_{F,\lambda+1} = \mathbf{y}_{F,\lambda} + h \sum_{i=1}^s b_i \mathbf{k}_{F,i}^\lambda, \quad (22)$$

$$\lambda = 0, 1, \dots, m-1,$$

$$\mathbf{k}_{F,i}^\lambda = \mathbf{f}_F \left(t_0 + \lambda h + c_i h, \mathbf{y}_{F,\lambda} + h \sum_{j=1}^s a_{ij} \mathbf{k}_{F,j}^\lambda, \tilde{\mathbf{Y}}_{L,i}^\lambda \right), \quad (23)$$

$$i = 1, 2, \dots, s,$$

with $\tilde{\mathbf{Y}}_{L,i}^\lambda \simeq \mathbf{y}_L(t_0 + \lambda h + c_i h)$.

The slowly varying vector is given by

$$\mathbf{y}_L(t_0 + H) \simeq \mathbf{y}_{L,1} = \mathbf{y}_{L,0} + H \sum_{i=1}^{\bar{s}} \bar{b}_i \mathbf{k}_{L,i}, \quad (24)$$

$$\mathbf{k}_{L,i} = \mathbf{f}_L \left(t_0 + \bar{c}_i H, \tilde{\mathbf{Y}}_{F,i}, \mathbf{y}_{L,0} + H \sum_{j=1}^{\bar{s}} \bar{a}_{ij} \mathbf{k}_{L,j} \right), \quad (25)$$

$$i = 1, 2, \dots, \bar{s},$$

with $\tilde{\mathbf{Y}}_{F,i} \simeq \mathbf{y}_F(t_0 + \bar{c}_i H)$.

From this definition we can attest that numerical coupling between fast and slow differential subsystems is achieved by $\tilde{Y}_{L,i}^\lambda$ and $\tilde{Y}_{F,i}$. Effective stratagems for computing these intermediate stage values are proposed in [8, 9].

3.3. Circuit-block partitioning strategy

We must be aware that circuit-block partition into fast and slow subsystems can possibly change during the time-step integration process. Hence, it will be very helpful if the simulator is capable of automatically detect the slow and fast variables in the circuit. This automatic classification can be achieved using embedded RK methods [2] and error estimates usually evaluated for step size control and stiffness detection [3]. Fast-slow partitioning strategies, step size control tools, number of microsteps within a macrostep, stiffness detection stratagems and many other technical details of MRK code implementation are thoroughly addressed in [8–10] and also in [11].

Finally, it must be pointed out that significant efficiency gains in computation speed for the simulation of several illustrative examples have been reported in the scientific literature, which demonstrate an effective reduction on the MRK computational cost in comparison to traditional SPICE simulation engines.

4. Multitime partitioned simulation engines

This section is devoted to briefly review some advanced circuit-block partitioning numerical simulation techniques operating in a multivariate time-domain framework. Section 4.1 introduces the multivariate formulation theory, in which the 1-D time is converted into a set of artificial time variables. Section 4.2 addresses some fundamental aspects of the numerical simulation algorithms. Finally, Section 4.3 describes some techniques for automatic circuit-block partition.

4.1. Multivariate formulation

The multivariate formulation is a useful stratagem that plays an important role in the EDA scientific community, especially in the RF and microwave areas. It was initially introduced in 1996 [12] as a sophisticated derivation of quasi-periodic HB, and it has been adopted by other researchers (e.g., [13–19]), which have demonstrated that it can be an efficient strategy when dealing with electronic circuits operating on widely distinct time scales. The success of multivariate formulation relies on the fact that voltages and currents containing components that evolve themselves at two, or more, widely separated rates of variation can be represented much more efficiently if we define them as functions of two, or more, time variables (artificial time scales). With this stratagem all signals (stimuli and responses) will be represented as multivariate functions, which will imply that dynamic behavior of the circuits will no longer be modeled by DAE, or ODE, systems formulated in the 1-D time. It will be described by partial differential systems.

In order to provide a simple and illustrative mathematical description of the multivariate formulation let us consider the classical example of a generic nonlinear RF circuit driven by an envelope modulated signal of the form

$$x(t) = e(t) \cos(\omega_C t + \phi(t)), \quad (26)$$

where $e(t)$ and $\phi(t)$ are, respectively, the amplitude and phase slowly varying baseband signals (envelope), modulating the $\cos(\omega_C t)$ fast-varying carrier. Computing the numerical solution of such circuit using conventional 1-D time-step integrators (RK methods or LMM methods) is very expensive. This is so because the response of the circuit has to be evaluated during a long time window defined by the slowly varying envelope, wherein the time-step size is severely restricted by the high frequency carrier. Taking into consideration the disparity between the baseband and the carrier time scales, and assuming that they are also uncorrelated, which is typically true, we are able to rewrite (26) as a bivariate function

$$\hat{x}(t_1, t_2) = e(t_1) \cos(\omega_C t_2 + \phi(t_1)), \quad (27)$$

where t_1 is the slow baseband artificial time scale and t_2 is the fast carrier artificial time scale. It must be noted that $\hat{x}(t_1, t_2)$ is a periodic function with respect to t_2 but not to t_1 , that is,

$$\hat{x}(t_1, t_2) = \hat{x}(t_1, t_2 + T_2), \quad T_2 = 2\pi/\omega_C. \quad (28)$$

This means that a generic 1-D $[0, t_{Final}]$ time interval will be mapped into a 2-D $[0, t_{Final}] \times [0, T_2]$ rectangular domain. It is easy to attest that, in general, the number of points required to represent $\hat{x}(t_1, t_2)$ in $[0, t_{Final}] \times [0, T_2]$ is much less than the number of points needed to represent $x(t)$ in $[0, t_{Final}]$. This is especially evident when the t_1 and t_2 time scales are widely separated [13, 14].

Let us now consider the DAE system of (1) modeling the dynamic behavior of a generic RF circuit excited by the envelope modulated signal of (26). Taking the abovementioned into account, we will adopt the following procedure for the vector-valued functions $x(t)$ and $y(t)$: t is replaced by t_1 in the slowly varying entities (envelope time scale) and t is replaced by t_2 in the fast-varying entities (RF carrier time scale). The application of this stratagem will turn the DAE system of (1) into the following partial differential system

$$p(\hat{y}(t_1, t_2)) + \frac{\partial q(\hat{y}(t_1, t_2))}{\partial t_1} + \frac{\partial q(\hat{y}(t_1, t_2))}{\partial t_2} = \hat{x}(t_1, t_2), \quad (29)$$

usually denoted as *multirate partial differential algebraic equations'* (MPDAE) system [13, 14]. It is easy to demonstrate that, if $\hat{x}(t_1, t_2)$ and $\hat{y}(t_1, t_2)$ satisfy (29), then the univariate forms $x(t) = \hat{x}(t, t)$ and $y(t) = \hat{y}(t, t)$ satisfy (1) [13]. Thus, $y(t)$ may be retrieved from its bivariate form $\hat{y}(t_1, t_2)$ by simply setting $t_1 = t_2 = t$, meaning that univariate solutions of (1) are available on diagonal lines $t_1 = t, t_2 = t$, along the bivariate solutions of (29). In truth, attending to the periodicity of the problem in the t_2 time dimension, the univariate version of the vector containing the circuit responses is obtained from its bivariate representation on the rectangular domain $[0, t_{Final}] \times [0, T_2]$ as

$$y(t) = \hat{y}(t, t \bmod T_2), \quad (30)$$

where $t \bmod T_2$ represents the remainder of division of t by T_2 .

The generalization of the bivariate strategy to a multidimensional problem with more than two time scales is straightforward. In fact, if the signals in the circuit present m separate rates of change, then m time scales will be used. In that case (29) assumes the generic form

$$p(\hat{\mathbf{y}}(t_1, t_2, \dots, t_m)) + \frac{\partial q(\hat{\mathbf{y}}(t_1, t_2, \dots, t_m))}{\partial t_1} + \dots + \frac{\partial q(\hat{\mathbf{y}}(t_1, t_2, \dots, t_m))}{\partial t_m} = \hat{\mathbf{x}}(t_1, t_2, \dots, t_m), \quad (31)$$

and the univariate solution, $\mathbf{y}(t)$, may be recovered from its multivariate form, $\hat{\mathbf{y}}(t_1, t_2, \dots, t_m)$, by setting $t_1 = t_2 = \dots = t_m = t$.

4.2. Partitioned algorithms for envelope following computation

The main advantage of the earlier described MPDAE approach is that it can result in significant improvements in simulation speed when compared to DAE-based alternatives [13–15]. However, by itself this approach does not perform any distinction between nodes or blocks in the circuit under analysis. In fact, in the first multivariate circuit simulation schemes initially introduced in [12], and then in [13], the same numerical algorithm was used to compute all the unknowns of the circuit. Only a few years later other advanced multivariate algorithms were proposed (e.g., [16–19]) in way to take into account possible circuit's heterogeneities. These algorithms are based on circuit-block partitioning stratagems defined within the multivariate frameworks. The most important ones regarding pure time-domain operations are briefly reviewed in the following.

As stated earlier, envelope modulated regimes are typical cases of practical interest. Computing responses to excitations of the form of (26) is a technique generally referred to as *envelope following*, which correspond to a combination of initial and periodic boundary conditions in the bivariate framework, leading to the following initial-boundary value problem

$$\begin{aligned} p(\hat{\mathbf{y}}(t_1, t_2)) + \frac{\partial q(\hat{\mathbf{y}}(t_1, t_2))}{\partial t_1} + \frac{\partial q(\hat{\mathbf{y}}(t_1, t_2))}{\partial t_2} &= \hat{\mathbf{x}}(t_1, t_2) \\ \hat{\mathbf{y}}(0, t_2) &= \mathbf{g}(t_2) \\ \hat{\mathbf{y}}(t_1, 0) &= \hat{\mathbf{y}}(t_1, T_2) \end{aligned} \quad (32)$$

defined on the rectangular domain $[0, t_{Final}] \times [0, T_2]$. $\mathbf{g}(\cdot)$ is a vector-valued initial-condition satisfying $\mathbf{g}(0) = \mathbf{g}(T_2) = \mathbf{y}(0)$, and $\hat{\mathbf{y}}(t_1, 0) = \hat{\mathbf{y}}(t_1, T_2)$ is the periodic boundary condition due to the periodicity of the problem in the t_2 fast time scale. In order to solve this initial-boundary value problem let us begin to consider the following semi-discretization of $[0, t_{Final}] \times [0, T_2]$ in the t_1 slow time dimension defined by

$$0 = t_{1,0} < t_{1,1} < \dots < t_{1,i-1} < t_{1,i} < \dots < t_{1,K_1} = t_{Final}, \quad h_{1,i} = t_{1,i} - t_{1,i-1}, \quad (33)$$

in which K_1 is the total number of steps in t_1 . Now, using a backward differentiation formula (BDF) [1] to approximate the t_1 derivatives of the MPDAE (for simplicity let us consider here the Gear-2 rule [1]), we obtain for each slow time instant $t_{1,i}$, from $i = 1$ to $i = K_1$, the periodic boundary value problem defined by

$$\begin{aligned} p(\widehat{\mathbf{y}}_i(t_2)) + \frac{3q(\widehat{\mathbf{y}}_i(t_2)) - 4q(\widehat{\mathbf{y}}_{i-1}(t_2)) + q(\widehat{\mathbf{y}}_{i-2}(t_2))}{2h_{1,i}} + \frac{dq(\widehat{\mathbf{y}}_i(t_2))}{dt_2} = \widehat{\mathbf{x}}(t_{1,i}, t_2), \\ \widehat{\mathbf{y}}_i(0) = \widehat{\mathbf{y}}_i(T_2), \end{aligned} \quad (34)$$

where $\widehat{\mathbf{y}}_i(t_2) \simeq \widehat{\mathbf{y}}(t_{1,i}, t_2)$. This means that, once $\widehat{\mathbf{y}}_{i-1}(t_2)$ is known, the solution on the next slow time instant, $\widehat{\mathbf{y}}_i(t_2)$, is obtained by solving (34). Hence, a set of K_1 boundary value problems have to be solved for obtaining the whole bivariate solution in the entire domain $[0, t_{Final}] \times [0, T_2]$. A very powerful technique has been proposed in the literature for solving each of the periodic boundary problems defined by (34) in an efficient way [16, 17]. This technique uses shooting-Newton based on MRK schemes. It splits the circuits into two distinct subsets according to the time evolution rates of their voltages and currents and performs time-step integration with different step lengths in each of the consecutive shooting iterations needed to solve (34). For that (34) is firstly divided into the following coupled fast-slow subsystems

$$\begin{aligned} p_F(\widehat{\mathbf{y}}_{F,i}(t_2), \widehat{\mathbf{y}}_{L,i}(t_2)) \\ + \frac{3q_F(\widehat{\mathbf{y}}_{F,i}(t_2), \widehat{\mathbf{y}}_{L,i}(t_2)) - 4q_F(\widehat{\mathbf{y}}_{F,i-1}(t_2), \widehat{\mathbf{y}}_{L,i-1}(t_2)) + q_F(\widehat{\mathbf{y}}_{F,i-2}(t_2), \widehat{\mathbf{y}}_{L,i-2}(t_2))}{2h_{1,i}} \\ + \frac{dq_A(\widehat{\mathbf{y}}_{F,i}(t_2), \widehat{\mathbf{y}}_{L,i}(t_2))}{dt_2} = \widehat{\mathbf{x}}(t_{1,i}, t_2), \\ p_L(\widehat{\mathbf{y}}_{F,i}(t_2), \widehat{\mathbf{y}}_{L,i}(t_2)) \\ + \frac{3q_L(\widehat{\mathbf{y}}_{F,i}(t_2), \widehat{\mathbf{y}}_{L,i}(t_2)) - 4q_L(\widehat{\mathbf{y}}_{F,i-1}(t_2), \widehat{\mathbf{y}}_{L,i-1}(t_2)) + q_L(\widehat{\mathbf{y}}_{F,i-2}(t_2), \widehat{\mathbf{y}}_{L,i-2}(t_2))}{2h_{1,i}} \\ + \frac{dq_L(\widehat{\mathbf{y}}_{F,i}(t_2), \widehat{\mathbf{y}}_{L,i}(t_2))}{dt_2} = \widehat{\mathbf{x}}(t_{1,i}, t_2), \end{aligned} \quad (35)$$

with

$$\widehat{\mathbf{y}}_i(t_2) = \begin{bmatrix} \widehat{\mathbf{y}}_{F,i}(t_2) \\ \widehat{\mathbf{y}}_{L,i}(t_2) \end{bmatrix}, \quad \widehat{\mathbf{y}}_{F,i}(t_2) \in \mathbb{R}^{n_F}, \quad \widehat{\mathbf{y}}_{L,i}(t_2) \in \mathbb{R}^{n_L}, \quad n_F + n_L = n. \quad (36)$$

$\widehat{\mathbf{y}}_{F,i}(t_2)$ is the vector containing the fast-varying circuit variables at the slow time instant $t_{1,i}$, and $\widehat{\mathbf{y}}_{L,i}(t_2)$ is the vector holding the slowly varying ones at the same slow time instant. The former will be integrated in t_2 with a small step size h_2 , while the later will be integrated with a much larger step size H_2 .

4.3. Circuit-block partitioning strategy

Since the partition of the electronic system in fast and slow subsystems may dynamically vary with time, it will be of great usefulness if the simulation algorithm is capable of automatically

detecting the fast-varying and the slowly varying signals. The approach adopted in [16, 17] for that purpose is briefly described in the following.

As mentioned earlier, each of the periodic boundary value problems defined by (34) is solved using shooting-Newton based on multirate time-step integrators (MRK schemes). Now, taking into account that shooting is an iterative technique, the key idea of this partitioning strategy is to use a uni-rate scheme on the first shooting iterations needed for each time line $t_{1,i} \times [0, T_2]$ of the $[0, t_{Final}] \times [0, T_2]$ rectangular domain. For that, it starts by considering $m = 1$ in the adopted MRK scheme, so that it degenerates into a standard uni-rate RK scheme. This means that all the unknowns (voltages or currents) will be integrated in t_2 with the same microstep h_2 . After that, the differential system of (34) is partitioned into (35) according to the variations in the t_2 derivatives of the unknowns. Each unknown that practically evidences no variations in its t_2 time derivatives for the entire time line $t_{1,i} \times [0, T_2]$, that is, each unknown that satisfies the condition

$$\left| \max(\text{slope}_j) - \min(\text{slope}_j) \right| < Tol, \tag{37}$$

where Tol is a small specified tolerance, will be treated as slow on the next shooting iterations. The remaining unknowns will be treated as fast. In a cheap scheme within a uniform grid each slope can be simply given by

$$\text{slope}_j = [\hat{y}_i(t_{2,j} + h_2) - \hat{y}_i(t_{2,j})]/h_2, \quad t_{2,j} \in \{0, h_2, 2h_2, \dots, T_2 - h_2\}. \tag{38}$$

All the subsequent shooting iterations on the time line $t_{1,i} \times [0, T_2]$ will be conducted in a multirate way using different step sizes. In a nonuniform grid, step sizes h_2 and $H_2 = m \cdot h_2$ may be chosen using any step-size control tool. In a uniform grid, they can be predefined or successively refined to achieve a desired accuracy.

The robustness of this partitioning strategy can be improved if more than one shooting iteration with the uni-rate scheme is considered for each time line $t_{1,i} \times [0, T_2]$. However, such tactic will obviously conduct to some extra computational cost, turning the overall algorithm to be a little less efficient.

Finally, to conclude this section, it must be highlighted that although efficiency gains are essentially dependent on the ratio between the number of slow and fast signals, significant reductions on the computational effort were reported in the scientific literature (e.g., [16, 17]) for the simulation of several illustrative examples of practical interest using the envelope following bivariate partitioned techniques.

5. Hybrid time-frequency partitioned techniques

This section is intended to provide a brief description of powerful circuit-block partitioning numerical simulation techniques operating in multivariate hybrid time-frequency frameworks. Section 5.1 introduces the fundamentals of multivariate hybrid envelope following. Section 5.2

addresses some basic details of the simulation algorithms and finally, Section 5.3 describes the approach for automatic circuit-block partition.

5.1. Multivariate hybrid envelope following

Let us once again consider the initial-boundary value problem of (32) characterizing the bivariate nature of an electronic circuit operating at two widely separated time scales. Let us also consider the semi-discretization of the $[0, t_{Final}] \times [0, T_2]$ rectangular domain, which has led to the set of periodic boundary value problems defined by (34). When waveforms are not excessively demanding on the quantity of harmonic components needed for an accurate frequency domain representation, we are able to use the HB method for efficiently computing the solution of (34). In such case we will obtain for each time line $t_{1,i} \times [0, T_2]$ the following HB system

$$\mathbf{P}(\widehat{\mathbf{Y}}(t_{1,i})) + \frac{3\mathbf{Q}(\widehat{\mathbf{Y}}(t_{1,i})) - 4\mathbf{Q}(\widehat{\mathbf{Y}}(t_{1,i-1})) + \mathbf{Q}(\widehat{\mathbf{Y}}(t_{1,i-2}))}{2h_{1,i}} + j\boldsymbol{\Omega}\mathbf{Q}(\widehat{\mathbf{Y}}(t_{1,i})) = \widehat{\mathbf{X}}(t_{1,i}), \quad (39)$$

in which $\widehat{\mathbf{X}}(t_{1,i})$ is a $(2K+1) \times n$ vector containing the Fourier coefficients of the stimuli (independent sources) and $\widehat{\mathbf{Y}}(t_{1,i})$ a $(2K+1) \times n$ vector containing the Fourier coefficients of the unknowns (node voltages and brunch currents' waveforms), at $t_1 = t_{1,i}$. K is the maximum harmonic order and n is the number of unknowns. In order to obtain the solution in the entire $[0, t_{Final}] \times [0, T_2]$ rectangular domain a total of K_1 HB systems have to be solved.

With this hybrid time-frequency approach we have an envelope following technique that handles the slow variations of the unknowns in the time domain and their fast variations in the frequency domain. This technique, usually denoted as multivariate *envelope transient harmonic balance* [20, 21], is able to exploit the existence of time-rate disparities in circuits operating under moderately nonlinear regimes of operation.

5.2. Partitioned algorithms for hybrid envelope following computation

The hybrid time-frequency envelope following technique presented earlier does not perform any distinction between nodes or blocks in the circuit. Although it has been widely used, especially by the RF and microwave community, only a few years ago other versions regarding circuit-block partition were proposed [18, 19] in way to take into account possible circuit's heterogeneities. The most important aspects of these partitioning techniques are briefly reviewed in the following.

Let us rewrite (39) as

$$\begin{aligned} \mathbf{F}(\widehat{\mathbf{Y}}(t_{1,i})) = \mathbf{P}(\widehat{\mathbf{Y}}(t_{1,i})) + \frac{3\mathbf{Q}(\widehat{\mathbf{Y}}(t_{1,i})) - 4\mathbf{Q}(\widehat{\mathbf{Y}}(t_{1,i-1})) + \mathbf{Q}(\widehat{\mathbf{Y}}(t_{1,i-2}))}{2h_{1,i}} \\ + j\boldsymbol{\Omega}\mathbf{Q}(\widehat{\mathbf{Y}}(t_{1,i})) - \widehat{\mathbf{X}}(t_{1,i}) = 0. \end{aligned} \quad (40)$$

The Newton iterative solver is usually utilized to solve (40), leading to

$$\mathbf{F}\left(\widehat{\mathbf{Y}}^{[r]}(t_{1,i})\right) + \left. \frac{d\mathbf{F}\left(\widehat{\mathbf{Y}}(t_{1,i})\right)}{d\widehat{\mathbf{Y}}(t_{1,i})} \right|_{\widehat{\mathbf{Y}}(t_{1,i})=\widehat{\mathbf{Y}}^{[r]}(t_{1,i})} \left[\widehat{\mathbf{Y}}^{[r+1]}(t_{1,i}) - \widehat{\mathbf{Y}}^{[r]}(t_{1,i}) \right] = 0, \quad (41)$$

which implies that for computing $\widehat{\mathbf{Y}}^{[r+1]}(t_{1,i})$ from the previous estimate $\widehat{\mathbf{Y}}^{[r]}(t_{1,i})$ we have to solve a linear system composed of $n \times (2K + 1)$ equations. The system of (41) involves the so-called Jacobian matrix of $\mathbf{F}\left(\widehat{\mathbf{Y}}(t_{1,i})\right)$, which has a block structure, consisting of an $n \times n$ matrix of square submatrices (blocks), each one with dimensions $(2K + 1) \times (2K + 1)$. The general block of row m and column l can be expressed as

$$\begin{aligned} \frac{d\mathbf{F}_m\left(\widehat{\mathbf{Y}}(t_{1,i})\right)}{d\widehat{\mathbf{Y}}_l(t_{1,i})} &= \frac{d\mathbf{P}_m\left(\widehat{\mathbf{Y}}(t_{1,i})\right)}{d\widehat{\mathbf{Y}}_l(t_{1,i})} + \frac{1}{h_{1,i}} \frac{d\mathbf{Q}_m\left(\widehat{\mathbf{Y}}(t_{1,i})\right)}{d\widehat{\mathbf{Y}}_l(t_{1,i})} \\ &+ j\Omega \frac{d\mathbf{Q}_m\left(\widehat{\mathbf{Y}}(t_{1,i})\right)}{d\widehat{\mathbf{Y}}_l(t_{1,i})}, \quad l, m = 1, 2, \dots, n, \end{aligned} \quad (42)$$

where $d\mathbf{P}_m\left(\widehat{\mathbf{Y}}(t_{1,i})\right)/d\widehat{\mathbf{Y}}_l(t_{1,i})$ and $d\mathbf{Q}_m\left(\widehat{\mathbf{Y}}(t_{1,i})\right)/d\widehat{\mathbf{Y}}_l(t_{1,i})$ denote the Toeplitz matrices [7] of the vectors containing the Fourier coefficients of $dp_m(\widehat{y}(t_{1,i}, t_2))/d\widehat{y}_l(t_{1,i}, t_2)$ and $dq_m(\widehat{y}(t_{1,i}, t_2))/d\widehat{y}_l(t_{1,i}, t_2)$, respectively.

A very powerful technique has been proposed in the literature [18, 19] for solving (41) in an efficient way. This technique takes into account that, in some cases (e.g., RF heterogeneous systems), there are parts of the circuits in which there are no fluctuations dictated by the fast carrier. As a consequence, bidimensional forms $\widehat{y}(t_1, t_2)$ of voltages and currents in those parts have no dependence on t_2 . This means that for each time line $t_{1,i} \times [0, T_2]$ each signal $\widehat{y}(t_{1,i}, t_2)$ has a constant value that can be represented as a Fourier series with a unique $k = 0$ coefficient. Thus, while quickly varying signals are represented in the frequency domain by a set of $(2K + 1)$ Fourier coefficients, signals in latent (slowly varying) blocks are represented by a single coefficient. Having this in mind it is straightforward to conclude that the size of vector $\widehat{\mathbf{Y}}(t_{1,i})$ will be significantly reduced, as well as the size (number of equations) of the HB system. Significant Jacobian matrix size reductions will also be achieved, since some of the submatrices will no longer be $(2K + 1) \times (2K + 1)$ matrices, but, instead, simple 1×1 scalar elements. For instance, let us suppose that the m th component of (40) is a slowly varying signal exclusively dependent on other slowly varying signals. Let us also suppose that the l th component is also a slowly varying signal. Then, the Jacobian matrix block of row m and column l will be a 1×1 block given by

$$\frac{d\mathbf{F}_m\left(\widehat{\mathbf{Y}}(t_{1,i})\right)}{d\widehat{\mathbf{Y}}_l(t_{1,i})} = \frac{dp_m(\widehat{y}(t_{1,i}))}{d\widehat{y}_l(t_{1,i})} + \frac{1}{h_{1,i}} \frac{dq_m(\widehat{y}(t_{1,i}))}{d\widehat{y}_l(t_{1,i})}. \quad (43)$$

It may be noted that this 1×1 scalar block can be seen as a particular case of the general $(2K + 1) \times (2K + 1)$ block of (42), if $K = 0$ is assumed as the maximum harmonic order. Actually, given that $df_m(\hat{y}(t_{1,i}, t_2))/d\hat{y}_l(t_{1,i}, t_2)$ is a constant function evidencing no fluctuations in the t_2 fast time scale, the last term of (42) will vanish and there will be no more necessity of converting the right-hand side terms of (43) into the frequency domain.

Since only fast-varying signals are converted into the frequency domain, this partitioned technique can be seen as a combination of multivariate envelope transient harmonic balance (by treating the fast-varying signals in a hybrid time-frequency framework) with a pure time marching simulation engine (by treating some of the signals in a pure time domain scheme). Thus, beyond the notorious significant vector and matrix size reductions above mentioned, there is another important advantage brought by this partitioned technique. For example, in complex heterogeneous RF systems strongly nonlinear regimes of operation are in general associated to digital or baseband blocks, whereas moderately nonlinear regimes are typical of RF blocks. With this partitioned technique signals in digital and baseband blocks are appropriately computed in the time domain, while signals in RF blocks are treated in the hybrid time-frequency framework.

5.3. Circuit-block partitioning strategy

Similar to what we have mentioned for the methods discussed in the previous sections, it will be of great utility if the simulator is able to automatically distinguish the fast-varying variables from the slowly varying ones. We now briefly review the approach addressed in [19] for that purpose, which splits the circuit into distinct blocks according to the time rates of change of their voltages and currents.

Let us consider the HB system of (39). As stated earlier, this system has to be solved with the iterative scheme of (41) for each artificial time line $t_{1,i} \times [0, T_2]$. The automatic partitioning strategy proposed in [19] consists in considering all the circuit's variables as fast on the first iteration of (41), that is, it consists in initially representing all the unknowns in the frequency domain as a set of $(2K + 1)$ Fourier coefficients. This way, the algorithm starts by computing a single iteration in (41) to evaluate $\hat{\mathbf{Y}}^{[1]}(t_{1,i})$

$$\hat{\mathbf{Y}}^{[1]}(t_{1,i}) = \left[\hat{\mathbf{Y}}_1^{[1]}(t_{1,i})^T, \hat{\mathbf{Y}}_2^{[1]}(t_{1,i})^T, \dots, \hat{\mathbf{Y}}_n^{[1]}(t_{1,i})^T \right]^T, \quad (44)$$

where each $\hat{\mathbf{Y}}_v(t_{1,i})$, $v = 1, \dots, n$, is a $(2K + 1) \times 1$ vector defined as

$$\hat{\mathbf{Y}}_v^{[1]}(t_{1,i}) = \left[Y_{v,-K}^{[1]}(t_{1,i}), \dots, Y_{v,0}^{[1]}(t_{1,i}), \dots, Y_{v,K}^{[1]}(t_{1,i}) \right]^T. \quad (45)$$

Each of the $\hat{\mathbf{Y}}_v^{[1]}(t_{1,i})$ is then inspected. If its Fourier coefficients of order $k \neq 0$ are practically null (their absolute values stay under a very small prescribed tolerance) it will be classified as slow. Otherwise it will be classified as fast. After this classification, which will temporarily split the system into fast and slow subsystems, the simulator considers that signals in slow subsystems can be represented by a single Fourier coefficient for the remaining iterations needed to compute the solution of (41).

Similar to what we have discussed for the method described in Section 4, the robustness of this partitioning strategy may be improved if more than one iteration in (41) is computed before the simulator decides which signals are slow and which signals are fast. The main drawback of such approach is the loss of some efficiency due to the extra computational effort required.

As a final point of this section, we would like to point out that significant gains in computation speed have been reported in the scientific literature (e.g., [18, 19]) for the simulation of several illustrative examples of practical interest using these hybrid time-frequency envelope following partitioned algorithms.

6. Conclusions

In this chapter, we have briefly reviewed some powerful numerical simulation techniques based on partitioned stratagems. Such techniques were especially designed to cope with the simulation challenges brought by emerging electronic technologies to the EDA community, as is the case of complex heterogeneous electronic systems composed of a combination of different kinds of circuit blocks (analog, mixed-signal and digital blocks, or even radio frequency blocks) containing node voltages and branch currents of very distinct formats and running on widely separated time scales. With these partitioned techniques signals within different blocks of the circuits are computed with distinct algorithms and/or step sizes. Considerable reductions on the computational cost have been registered in several experiments published in the scientific literature (in comparison to previously recognized techniques) without compromising the accuracy of the results.

Acknowledgements

This work is funded by National Funds through FCT—Fundação para a Ciência e Tecnologia, under the project UID/EEA/50008/2013.

Author details

Jorge dos Santos Freitas de Oliveira^{1,2*}

*Address all correspondence to: jorge.oliveira@ipleiria.pt

1 School of Technology and Management, Polytechnic Institute of Leiria, Leiria, Portugal

2 Institute of Telecommunications, University of Aveiro, Aveiro, Portugal

References

- [1] Lambert J. Numerical Methods for Ordinary Differential Systems: The Initial Value Problem. 1st ed. West Sussex: Wiley; 1991

- [2] Hairer E, Nørsett S, Wanner G. Solving Ordinary Differential Equations I: Nonstiff Problems. 1st ed. Berlin: Springer-Verlag; 1987. DOI: 10.1007/978-3-662-12607-3
- [3] Hairer E, Wanner G. Solving Ordinary Differential Equations II: Stiff and Differential Algebraic Problems. 2nd ed. Berlin: Springer-Verlag; 1996. DOI: 10.1007/978-3-642-05221-7
- [4] Kundert K, White J, Sangiovanni-Vincentelli A. Steady-State Methods for Simulating Analog and Microwave Circuits. Norwell, MA: Spinger; 1990
- [5] Rodrigues P. Computer-Aided Analysis of Nonlinear Microwave Circuits. Norwood, MA: Artech House; 1998
- [6] Maas S. Nonlinear Microwave and RF Circuits. 2nd ed. Norwood, MA: Artech House; 2003
- [7] Pedro J, Carvalho N. Intermodulation Distortion in Microwave and Wireless Circuits. 1st ed. Norwood MA: Artech House; 2003
- [8] Kværnø A. Stability of multirate Runge-Kutta schemes. *International Journal of Differential Equations and Applications*. 2000;1A:97-105
- [9] Günther M, Kværnø A, Rentrop P. Multirate partitioned Runge-Kutta methods. *BIT*. 2001; 41(3):504-514. DOI: 10.1023/A:1021967112503
- [10] Bartel A, Günther M, Kværnø A. Multirate methods in electrical circuit simulation. *Progress in Industrial Mathematics*, Springer. 2002;1:258-265. DOI: 10.1007/978-3-662-04784-2_35
- [11] Oliveira J, Araújo A. Envelope transient simulation of nonlinear electronic circuits using multirate Runge-Kutta algorithms. *WSEAS Transactions on Electronics*. 2006;3(2):77-84
- [12] Brachtendorf H, Welsch G, Laur R, Bunse-Gerstner A. Numerical steady-state analysis of electronic circuits driven by multi-tone signals. *Electrical Engineering (Springer-Verlag)*. 1996;79(2):103-112. DOI: 10.1007/BF01232919
- [13] Roychowdhury J. Analyzing circuits with widely separated time scales using numerical PDE methods. *IEEE Transactions on Circuits and Systems*. 2001;48(5):578-594. DOI: 10.1109/81.922460
- [14] Mei T, Roychowdhury J, Coffey T, Hutchinson S, Day D. Robust stable time-domain methods for solving MPDEs of fast/slow systems. *IEEE Trans. on Computer-Aided Design of Integrated Circuits and Systems*. 2005;24(2):226-239. DOI: 10.1109/TCAD.2004.841073
- [15] Zhu L, Christoffersen C. Adaptive harmonic balance analysis of oscillators using multiple time scales. In: 3rd International IEEE Northeast Workshop on Circuits and Systems; Québec City. 2005. pp. 187-190
- [16] Oliveira J, Pedro J. An efficient time-domain simulation method for multirate RF nonlinear circuits. *IEEE Transactions on Microwave Theory and Techniques*. 2007;55(11): 2384-2392. DOI: 10.1109/TMTT.2007.908679
- [17] Oliveira J, Pedro J. A multiple-line double multirate shooting technique for the simulation of heterogeneous RF circuits. *IEEE Transactions on Microwave Theory and Techniques*. 2009;57(2):421-429. DOI: 10.1109/TMTT.2008.2011228

- [18] Oliveira J, Pedro J. A new mixed time-frequency simulation method for nonlinear heterogeneous multirate RF circuits. In: Proceeding of the IEEE MTT-S International Microwave Symposium (MTT '10); May 2010; Anaheim, CA. pp. 548-551
- [19] Oliveira J, Pedro J. Efficient RF circuit simulation using an innovative mixed time-frequency method. IEEE Transactions on Microwave Theory and Techniques. 2011;59(4): 827-836. DOI: 10.1109/TMTT.2010.2095035
- [20] Rizzoli V, Neri A, Matri F. A modulation-oriented piecewise harmonic-balance technique suitable for transient analysis and digitally modulated analysis. In: Proc. 26th European Microwave Conference; Jun. 1996; Prague. 1996. pp. 546-550
- [21] Sharrit D. Method for simulating a circuit. U.S. Patent 5588142, December 24, 1996

Numerical Simulation of Fission Product Behavior Inside the Reactor Containment Building Using MATLAB

Khurram Mehboob and
Mohammad Subian Aljohani

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.70706>

Abstract

The aim of this work is to carry out the numerical simulation of fission product (FP) behavior inside the reactor building under loss of coolant accident (LOCA) using MATLAB. For this purpose, a kinetic model has been developed and implemented in MATLAB to study the behavior of in-containment FPs during postulated LOCA for typical 1000 MW pressurized water reactor (PWR). A continuous release of the FPs from the reactor pressure vessel (RPV) has been implemented with coolant retention. The in-containment FP behavior is influenced by containment atmosphere and containment safety systems. The sensitivity analysis and removal rate of airborne isotopes with the containment spray system have been studied for various spray activation time, spray failure time, droplet size and spray pH value. The droplet size and pH value of the spray system effectively remove the airborne isotopes. The alkaline (sodium thiosulfate, $\text{Na}_2\text{S}_2\text{O}_3$) spray solution and spray with pH 9.5 have similar scrubbing properties for iodine. However, the removal rate from the containment spray system has been found an approximately inverse square of droplet diameter ($1/d^2$).

Keywords: PWR, fission products, MATLAB, sodium thiosulfate, droplet diameter

1. Introduction

The nuclear reactor systems are sufficiently complex that there could be the possibility of an accident followed by the release of fission product (FPs). Such a release could require multiple failures of safety systems and barriers. In the case of a break in the hot/cold leg in a pressurized water reactor (PWR), coolant and energies are first released from the reactor coolant system to the containment through the break. The FP also released along with the coolant. This type of accident usually occurs in the high-pressure cold leg. The worst

condition of such an uncontrolled break is the guillotine type of break. In such type of accident, the envelope of primary systems is breached [1]. If such an accident is not controlled by safety systems, then such accidents may transform into the severe accident.

In severe accidents, FP is released during the progression of accidents [2]. Owing to the strong influence of thermal hydraulics on FP release and transportation, FP release and transport mechanism is very complicated and complex. The FP behavior inside the containment is the fundamental of the source term. The source term results are the outputs of level 2 PSA [3], which are necessary for radiological assessments and consequences. The dominant FPs that constitute in hazardous effects can be categorized as noble gasses (Xe, Kr), volatile (I, Cs, Te), semi-volatile (Ru, Ag, Ba, Sr, Tc, Rh) and nonvolatile (Nb, Zr, Y, Pd, La, Mo, Tc, Nd, Ce) ([4, 5]). The aerosols are ^{129}Te , ^{127}Te , ^{105}Rh , ^{103}Ru , ^{105}Ru , ^{137}Cs , ^{138}Cs , ^{89}Sr , ^{90}Sr and ^{140}Ba . These isotopes release in the particulate form, and going through agglomeration and nucleation process, they form aerosols [6]. However, iodine may transform into volatile species and possess a complex chemistry [7]. The common organic form of iodine is available in chemical forms as CH_3I , CsI and HI [8]. The behavior of FP is highly influenced by the in-containment atmosphere, heat loads, containment pressure and steam generation rate. The containment is installed with the spray system and cooling fans to prevent the early over-pressurization due to the heat load. The containment spray system is significant in enhancing the early depletion of radionuclides during early in-vessel release phase from the containment atmosphere. The spray system is automatically activated, as an emergency designed device to prevent containment integrity [9].

The FP release from a nuclear power plant (NPP) is known as a key factor affecting both the design of safety equipment and safety evaluation, including safety and risk assessment [10]. Experimental research on FP release behavior was conducted by many investigators [11–13]. Many experiments had played a significant role in understanding the behavior of aerosols, FPs, iodine chemistry, and transportation under accident situations [13–17]. The Phébus-FP project [18] was the most impressive program initiated to study the behavior of FP. The main objectives of this project were (1) to minimize the uncertainty in source term evaluation, (2) to study the FPs, structural and control rod material release transportation and deposition from the degraded core through coolant, and (3) and behavior of FP inside the containment building [19, 20]. Meanwhile, several analytical and computational codes were developed. ASTEC is one of the most popular codes used to study the behavior of FPs in severe accident conditions [21]. MELCOR along with MACCS can be used to assess FP release and assessment of radiological consequence [22]. MAAP is the most popular tool to calculate severe accident source term, and its quick calculation is its prime character. Therefore, MAAP code is widely used in the nuclear industry [23].

Moreover, the numerical simulation of FP activity has been carried out by several researchers. [24] have developed an analytical model FIPRAP “FP Release Analysis Program” for the numerical simulation of FPs released from the fuel. The FIPRAP code can estimate the volatile FPs released from the nuclear fuel under changing irradiation conditions with the incorporation of all physical phenomena and fulfill the requirements of fuel designing, performance, degradation and source term estimation codes. Lewis et al. [25] have presented a review of FPs release modeling in support of fuel failure monitoring analysis for the characterization and allocation of defected fuel. A generalized model for FP transport in the fuel-to-sheath gap was

given by [26]. Koo et al. [27] have proposed a model describing pallet oxidation and bubble formation at grain boundaries, their interlinkages, and release into exposed surfaces. Avano [28] described a good model for the description release of FPs from the porous ceramic fuel, its leakage from cladding and mixing with the primary coolant. Tucker and white [29] have proposed an analytical model for the estimation of FPs from ceramic UO₂ fuel. In this model, the PF leakage probabilities from the fuel interior through grain are figured out. These probabilities strongly depend on the interconnectivity of pores in the ceramic fuel. Awan et al. [30] have also carried out the numerical simulation of FP activity in the reactor primary coolant. The proposed developed model is hybrid and analyzes the static and dynamic FP activity in the primary coolant of the reactor [31].

The goal of this chapter is to carry out the numerical simulation of FP behavior inside the reactor containment building under LOCA using MATLAB. The calculation process of iodine and other FPs is shown in **Figure 1**. A semi-kinetic model has been developed and implemented in MATLAB to carry out the sensitivity analysis of FPs during postulated LOCA for typical 1000 MW PWR. The kinetic model is presented in section II, which contains the deterministic as well as the kinetic approach. The deterministic computational methodology and computational steps flow chart are described in Section III. Next, the flow chart of model and implementation of model in MATLAB are described in Section IV. The examples and outcomes of the simulation results are presented in Section V. Finally, Section VI is the conclusion.

2. In-containment fission product release model

Figure 1 shows the process of release of FPs from fuel to cladding, cladding to coolant and then to the containment. In this work, a 1000-MW pressurized water reactor (PWR) has been considered with the design specification as shown in **Table 1**. The PWR system along with the containment system is shown in **Figure 2**. We have developed a real-time kinetic model to simulate the FP behavior inside the containment. The analytical model is a set of coupled ordinary differential equations (ODEs). The FP activity inside the reactor containment building and on the surfaces and walls of the containment is governed by the following sets of ODEs [8, 32, 33].

$$\frac{dm_{v,i}(t)}{dt} = -\lambda_i m_{v,i}(t) - u_{t,i} \frac{S}{V} m_{v,i}(t) - \alpha \frac{F}{V} m_{v,i}(t) - R_{res,i} \frac{\eta_{rc}}{V} m_{v,i}(t) - \frac{L_r}{V} m_{v,i}(t) + r_i \frac{S}{V} m_{s,i}(t) + P_i(t) \quad (1)$$

where

$$\alpha = \begin{cases} H\eta_i & \text{Iodine} \\ \frac{3hEa}{2d} & \text{other FPs} \end{cases} \quad (2)$$

$$\frac{dm_s(t)}{dt} = v_i m_v(t) - r m_s(t) \quad (3)$$

where i indicates the isotope, whereas V and S indicate the volumetric and surface activities of ith isotope. The puff release of FP is $m_v(t) = f_x \times f_f \times f_p \times f_c \times A_c/V \text{ g m}^{-3}$. The values of various parameters used in these simulations are listed in **Table 2**.

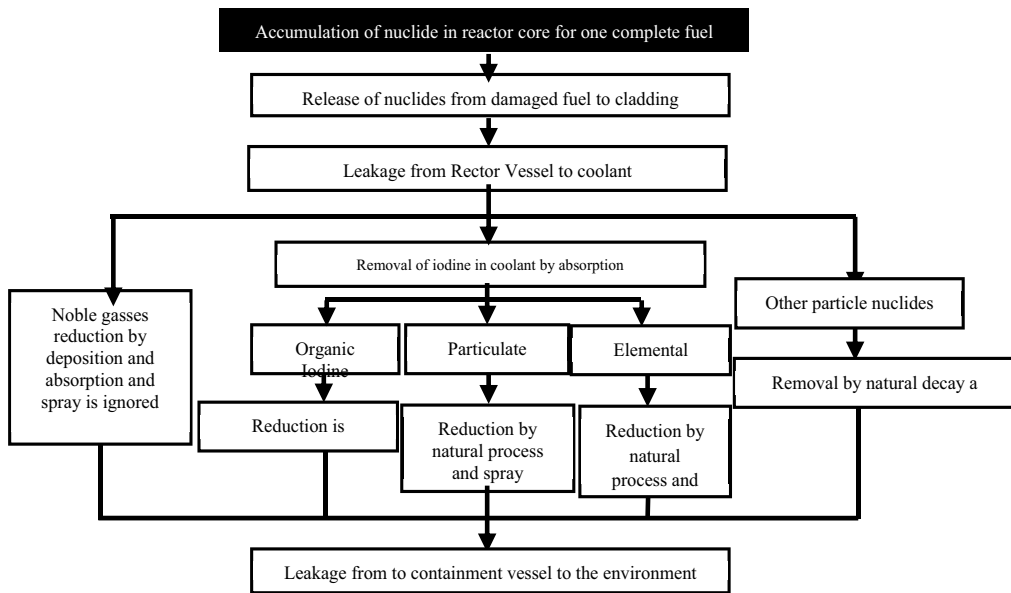


Figure 1. The calculation process of the FP behavior inside the reactor building.

Parameters	Value
Reactor	PWR
Fuel type	UO ₂
Average fuel enrichment wt%	2.4%
Specific power (MWth/kg U)	33.3
Power density (MWth/m ³)	66.6
System pressure (MPa)	15.166
System pressure (MPa)	14.96
Coolant flow (kg/s)	17387.7
Core height (m)	12.41
Core active region (m)	3.65
Core diameter (m ²)	3.81
Fuel assemblies	177
Control rod assemblies	69
Cladding material	Zircaloy
Fuel rod outer diameter (cm)	1.092
Rod pitch (cm)	1.443
Fuel assembly matrix	15 × 15
Coolant inlet temperature (k)	564.81

Parameters	Value
Coolant outlet temperature (k)	592.98
Control rods	1104
Control rod material	Ag (80%)-In (15%)-Cd (5%)

Table 1. Design parameters of typical 1000 MW reactor [34, 35].

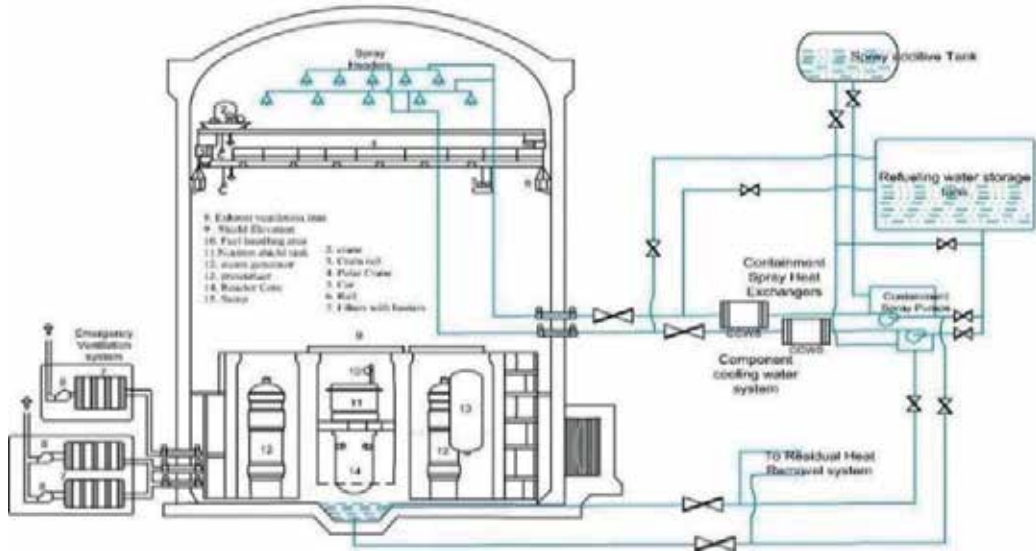


Figure 2. A schematic diagram of a typical PWR system with the containment spray system.

Parameters	Symbol	Value
Containment free volume	V (m^3)	57,600
Containment free surface	S (m^2)	~34,374
Leakage rate	L_r (m^3/s)	14.15
Core damage fraction	f_c (%)	35%
Fuel release fraction	f_f (%)	9.0×10^{-1}
Water release fraction	f_p (%)	$3.00 \times 10^{-1} - 1.0 \times 10^2$
Recirculation rate	R_{res} (m^3/s)	1-5
Recirculation filtration efficiency	η_{res} (%)	10-90%
Exhaust filter efficiency	η_{ex} (%)	90-98
Fraction of immediately released radioisotopes	f_x (%)	2.0×10^{-1}
Mixing rate	w_x (s^{-1})	0.1-1.0
Spray flow rate	F (m^3/s)	0.1-1.0

Parameters	Symbol	Value
Droplet size	d (micron)	100–1000
Deposition velocity (ud)	(m/s)	5.5×10^{-4}
Resuspension rate	s^{-1}	$\leq 2.3 \times 10^{-6}$

Table 2. Important parameters used for simulation [36].

2.1. Kinetic source of fission product

The last term in Eq. (1) is the source of FP from the reactor pressure vessel. The kinetic source is modeled as [37].

$$P(t) = (1 - f_x) A_c f_f f_p f_c \frac{K}{V} e^{(-w_x t)} \quad (4)$$

$$K = \frac{w_x \times (w_x / T)}{w_x - w_x / T} \quad (5)$$

The $(1 - f_x) \exp.(-w_x t)$ is the airborne FP activity released along with the coolant with mixing rate w_x . Where K is the normalization constant and expressed as follows. The overall radioactive mass inventory, including kinetic and static parts, is depicted in Eq. (6).

$$A_c = f_x A_c + (1 - f_x) A_c B \int_0^T e^{-w_x t} dt \quad (6)$$

2.2. Fission product removal with spray

The removal of iodine and aerosols from the containment with the spray system can be expressed as depicted in Eqs. (7) and (8), where m_{ri} and m_{ra} are the removal rates of iodine and aerosols, respectively.

$$\frac{dm_{ri,i}(t)}{dt} = P_i(t) - \frac{H\eta_i F}{V} m_{v,i}(t) \quad (7)$$

$$\frac{dm_{ra,i}(t)}{dt} = P_i(t) - \frac{3hFEa}{2dV} m_{v,i}(t) \quad (8)$$

where

$$\eta_i = 1 - e^{-6(K_G \times t_d / d \times (H + \kappa_G / \kappa_L))} \quad (9)$$

and

$$K_G = \frac{D_L}{d} \{2.0 + 0.60 \times \text{Re}^{0.5} \times \text{Sc}^{0.33}\} \quad (10)$$

$$K_L = \frac{2\pi^2 D_L}{3d} \quad (11)$$

$$D_L = \frac{(7.4 \times 10^{-8}) \times \sqrt{(xM_1)} \times T}{\mu_1 v^{0.6}} \quad (12)$$

The values of these parameters in Eqs. (9)–(12) are listed in **Table 3**.

Parameters	Symbols	Values
Partition coefficient	H	200 (for pH 5.0), 5000 (for pH 9.5) and 10,000 ($\text{Na}_2\text{S}_2\text{O}_3$)
Reynolds number	Re	1.29
Schmitt number	Sc	1.742
Molar weights of solvent	Ml (g/mole)	18.01528
Temperature	T (K)	$80 + 273.15$
Molecular volume of I_2	v (cm^3/g)	71.5
Viscosity	μl (centipoise)	0.35
Spray flow rate	F (m^3/sec)	0.35
Degree of solvent	x	2.6 for H_2O

Table 3. Numerical data for spray removal term ([36, 38]).

3. Deterministic computational methodology

Several steps are involved in the simulation of FP behavior inside the reactor building starting from the generation of FP in fuel along with the fuel burn-up. Leakage of FP into the coolant and then from the coolant to containment along with the leakage of coolant. The computational steps are listed in **Figure 3**. A two-stage methodology has

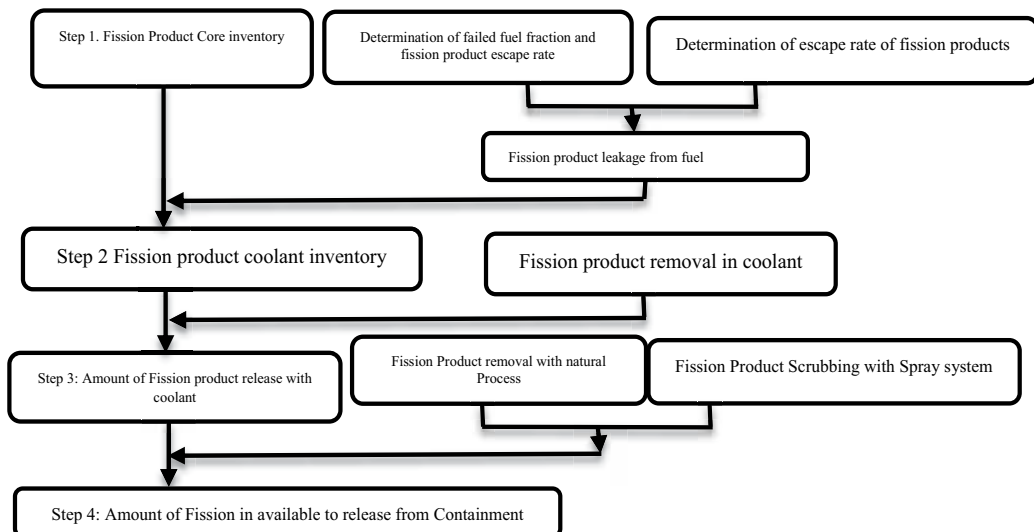


Figure 3. Flow chart of incontinent FP source term estimation.

been adopted: (1) evaluation of activity in the core just before the accident and (2) kinetic quantification of airborne activity under confined conditions. The core activity has been evaluated at for one complete fuel cycle to get maximum core activity. The behavior of airborne FP activity has been quantified for loss of the coolant accident (LOCA) under NUREG-1465 [8] and regulatory guide 1.183 [32] assumptions. The developed model uses subroutine functions containing coupled ODEs and Runge–Kutta (RK) method. The ODEs (Eqs. (1)–(12)) are implemented in MATLAB. The system of ODEs (Eqs. (1), (3), (7), (8)) is coupled and solved numerically using the Runge–Kutta (RK) method in this program.

The RK numerical provides efficient time-domain solution, yielding static as well as dynamic values of FPAs corresponding to about 84 different dominant FPs. The computational cycle starts with the initialization of the variables with $t = 0$. In the time loop, the values of FPAs inside the containment building are calculated using RK scheme for each next time step. The program allows performing these calculations for spray system operation.

4. Implementation in MATLAB

The above equations can be implemented in MATLAB. The flow chart of the MATLAB program is shown in **Figure 4**. In the first step, the physical constant and parameters are defined, and the time array and droplet size are determined by the user.

```
function PWR_Fission_Product
% MATLAB Program for In-containment Fission product program by Khurram Mehboob
% Date : 08-07-2017
%=====
clear; clc; clear all;
%=====
Global Hi Lr V S vd dec r Rr neu EI h Klcm Kgcm d Ea fr H y00 Q y t I Ac D Core_I
Cont_A QQ fx fc B wx YY Sorc wx1
tn = input('Enter end time = tn = '); h = input('Enter stepsize = h = '); t = (0:h:tn); % time array
for d1=100: 100: 1000; % particle diameter (microns)
%=====Control Variables=====
d = d1*1e-4; % particle diameter (cm)
k=d1/100; % Droplet control Factors for printing
fx = 0.20; % activity immediately available in the containment air
fc = 0.35; % core damage fraction.
H =10000; % partition coefficient for iodine
Rr = 4.719; % Recirculation flow rate
Lr = 14.15; % leakage rate
wx = 0.01; % mixing rate
```

In the second step, the fixed variables are loaded from an input text file. The input text file contains the output data from the ORIGEN2.2 code that contains data for 84 different FPs.

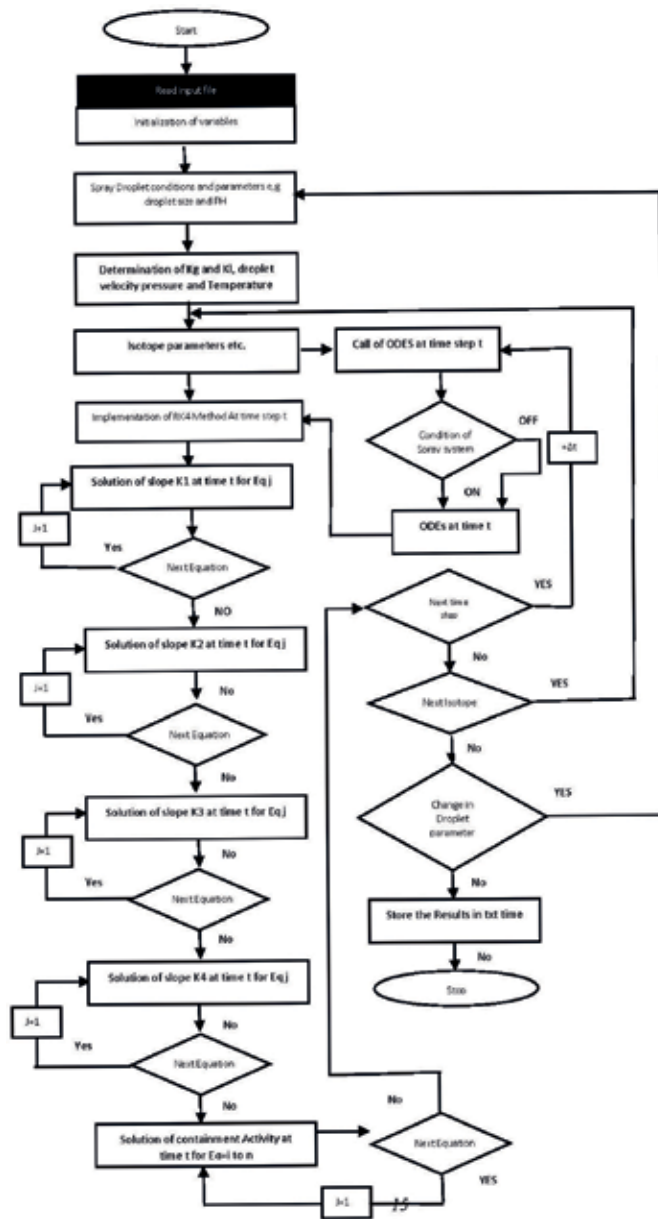


Figure 4. Flow diagram of computer program.

```
load 'input.txt'
%=====Fixed variables=====
V = input2(1,1);           % free volume of the containment
S = input2(2,1);           % free surface of the containment
No_iso = input2(5,1);      % no of isotopes
Hi = input2(6,1);          % height of spray system 40.0 m
fr = input2(10,1);         % Recirculation filtration efficiency
```

```

wx1 = wx/10.0;           % mixing rate of fission products in coolant
K = (wx*wx1)/(wx-wx1);  % normalization constant
%=====
for i = 1:No_iso;       % loop to read isotope data
Ac = input(i+10,4);    % fission product activity in the core.
    for j = 6:7
        Ac = Ac*input(i+10,j);  % Ac*ff*fw
    end
dec = input2(i+10,2);  % decay constant
vd = input2(i+10,08); % deposition velocity
r = input2(i+10,09);  % resuspension rate
neu = input2(i+10,10); % heap filter efficiency

```

In the third step, the fixed variables for Eqs. (9)–(12) are solved for droplet by using the data listed in **Table 3**. The values of parameters (T , x , M_l , u_l , Sc , Re) are the constant variable for the static containment atmosphere.

```

T = 80+273.15;          % Temperature (K)
x = 2.6;                % H2O
Ml = 18.01528 ;        % molar weight of solvent (g/mole )
v = 71.5;               % molar volume of diffusing substance (l2)
Sc = 1.742;             % Schmidt number.
Re = 1.95;              % Reynold number.
c1 = T-281.615; c2 = (T-281.615)^2; c3 = 8078.4+ c2; c4 = c3^0.5; C = 2.1482*(c1+c4)-120;
ul = (100/C)*2.41908832931* 14.8816390000001; %conversion of Centipoise to g/cm.s

```

Next, the fixed variables are used to determine the dynamics variables (DL , D_{mix} , KG , KL). The Eqs. (9)–(12) are solved using the static data calculated above. Also, the kinetic release of FP is $Q(t)$ that is determined as the function of time.

```

DI2= (((7.4)*(10^-8))*((x*Ml)^0.5)*T)/(ul*(v^0.6)) ; %diffusion coefficient if I2 (cm2/s)
DIMix = 0.00035*0.258064; %diffusion coeffi of I2 in steam ( cm2/s)
Kgcmm = (DIMix/d)*(2.0+0.6*(Re^(1/2))*(Sc^(1/3))); % liquid phase trans coefficient ( cm/s )
Klcm = ((2/3)*(3.1416^2)*(DI2))/(d); % gas phase trans coefficient (cm/s)
vt = ((Re*2387.4)/(d1))* 0.508; % terminal velocity (cm/s)
te = (40*100)/(2*vt); % exposure time (s)
A = Kg*(te); BB = d*(H+(Kg/Kl));
EI = (1 - exp(-6*(A/BB))); % FP removal rate (s-1)
SorC= Ac*fc/V;
Q = ((1-fx)*B*SorC*(exp(-wx1*t)- exp(-wx*t))); % kinetic source

```

Next, the initial conditions for airborne and surface activities are implemented, for example, $m_v(t) = f_x \times f_f \times f_p \times f_c \times A_c/V \text{ g.m}^{-3}$ and $m_s(t) = 0.0$. The kinetic and static parameter values are implemented in coupled equation written as subfunction diffeq, and Rang–Kutta fourth-order method is implemented by calling odeRK4 subroutine function.

```
%Initial Condition
%=====
y00 = fx*fc*Ac/V;    y0 = [y00, 0, 0];
%=====
[t,y] = odeRK4(@diffeq, tn, h, y0);
%=====
YY(:,i)= y(:,1) ;      % YY=y(length(y),:);
I(i)= input2(i+10,1);  % Atoms number
D(i)= input2(i+10,2);  % decay constant
Core_I(i)= input2(i+10,4); % Activity in the core
Cont_A(i)= y00;        % immediate released activity in containment
end
```

The subroutine function containing Eqs. (1), (3), (4) and (7) or (8) is depicted below. The condition for containment spray is controlled in subroutine function here.

```
function dy = diffeq(t,y)
global Lr V S vd dec r F EI H Rr fr fx Ac fc wx B wx1 Q Sorc
if t <=700
    F = 0.0; % input2(8,1); % spray flow rate (0.1-2.0 m3/s)(950 m3/h=0.264m3/s)
else
    F= 0.35;
end
Q = ((1-fx)*B*Sorc*(exp(-wx1*t)- exp(-wx*t)));
dy(1)= -dec*y(1)-(Lr/V)*y(1)-vd*(S/V)*y(1)+r*(S/V)*y(2)-fr*(Rr/V)*y(1)-((F*H*EI)/V)*y(3)+((1-fx)
*B*Ac*(fc/V)*((exp(-wx1*t))-exp(-wx*t)));
dy(2)= vd*y(1)- r*y(2);
dy(3)= Q - ((EI*H*F)/V)*y(3);
end
```

The Range–Kutta fourth-order method is implemented by calling the odeRK4 subroutine function. The function is capable of solving N number of coupled ODEs at separate time steps. The implementation if RK4 method is shown below.

```
% implementation of Range kutta 4thorder method.
function [t,y] = odeRK4(diffeq,tn ,h, y0)
t = (0:h:tn);          % time vector with spacing h
nt = length(t);        % no. of elements in vector t
neq = length(y0);
y = zeros(nt, neq);    % prelocation of y for speed
y(1,:) = y0(:);
h2=h/2;  h3=h/3; h6=h/6;
k1 = zeros(neq,1); k2= k1; k3 = k1; k4= k1; ytemp = k1;
for j=2:nt
```

```

told = t(j-1); yold = y(j-1,:);
k1= feval(diffeq,told, yold);
for n= 1:neq
    ytemp(n) = yold(n) + h2*k1(n);
end
k2= feval(diffeq,told+ h2, ytemp);
for n= 1:neq
    ytemp(n) = yold(n) + h2*k2(n);
end
k3= feval(diffeq,told +h2, ytemp);
for n= 1:neq
    ytemp(n) = yold(n) + h*k3(n);
end
k4= feval(diffeq,told+h, ytemp);
for n= 1:neq
    y(j,n)= yold(n)+h6*(k1(n)+k4(n))+h3*(k2(n)+k3(n));
end
end
end
end

```

5. Some results from numerical simulation

5.1. Volumetric fission product inventory

The core inventory for typical 1000 MW PWR has been evaluated by ORIGEN 2.2 code which is used by our model as a subroutine. A 35% core damage has been considered and 20% (f_x) as the puff release. While the rest of radioactive mass release along with coolant with mixing rate $w_x = 0.01 \text{ s}^{-1}$. The FP release inside the containment building as a function of time is depicted in **Figure 5**. The volumetric radioactive mass found to increase during

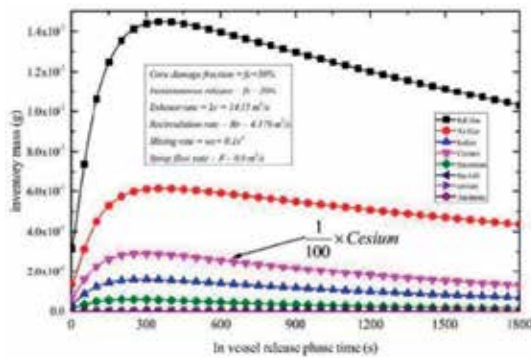


Figure 5. In-containment FP inventory during in-vessel release phase with mixing rate $w_x = 0.01 \text{ s}^{-1}$.

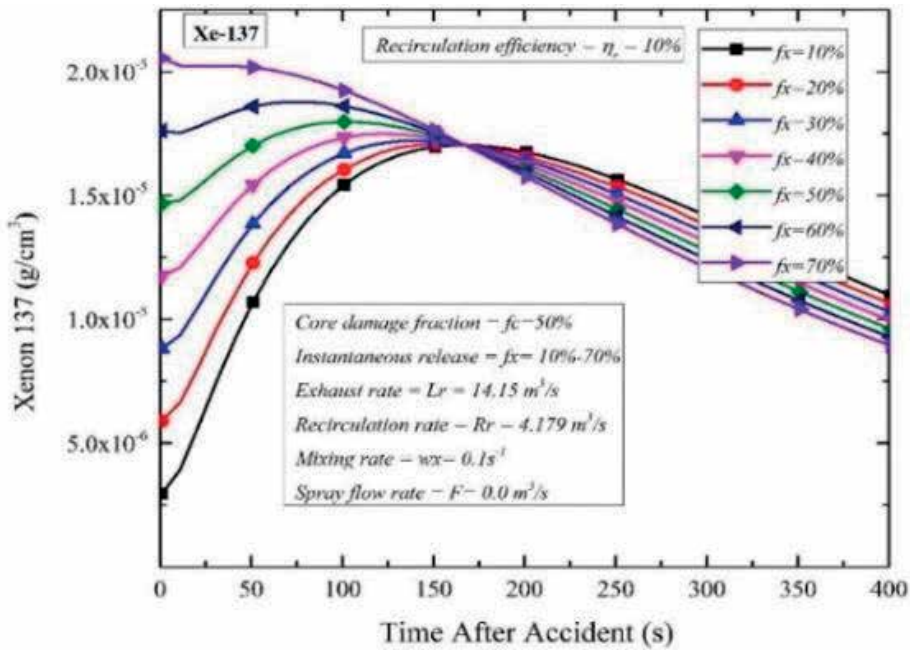


Figure 6. ^{137}Xe activity (g/cm^3) as function of time (s) for various values of f_x and $w_x = 0.01/\text{s}$.

first 300 s then starts decreasing with a constant rate. The cesium found to be dominant with 100 times higher than the other radioactive released masses. The Krypton gas is found to be 15% higher in magnitude with Xenon gas. However, the other isotopes show the similar behavior but with less in magnitude.

5.2. Puff release (f_x) effect on in-containment FPA

The LOCA is due to the uncontrolled leakage from coolant piping (hot leg or cold leg). The coolant burst release generates the immediate escape of radioiodine into containment with the rapture. The volumetric activity of ^{137}Xe inside the containment has been simulated for various values of instantaneous burst release (i.e., $f_x = 10\text{--}70\%$) of total activity inside the core. The simulation results are depicted in Figure 6. The results indicated that with the higher percentage of instantaneous release ($f_x = 50\text{--}70\%$), the activity in containment slightly increased and then decreased linearly.

However, a less fraction of burst release ($f_x = 10\text{--}30\%$ of total activity), the activity inside the containment first increases and then starts decreasing after approaching to the maximum value. As the value of f_x decreases, the peak shifts toward higher timescale. The peak becomes more prominent with small values of f_x . This happens due to competition between f_x term in the initial condition and $(1-f_x)(\exp -w_x)$ term in source term (Eq. (4)). The behavior of ^{137}Xe for various values of instantaneous release (f_x) with mixing rate $w_x = 0.01/\text{s}$ explains the clearer picture of airborne Xenon (Figure 6).

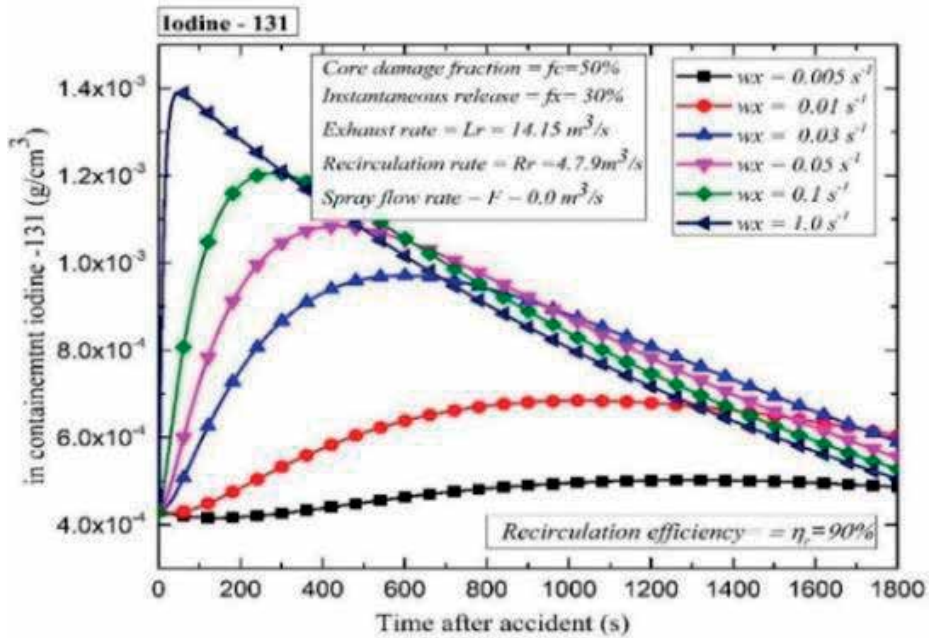


Figure 7. ^{131}I activity (g/cm^3) for ^{131}I as a function of time (s) for various values of w_x .

5.3. Mixing rate (w_x) effect on in-containment FPA

The delayed core released fraction $(1-f_x)(\exp -w_x t)$ for various values of mixing rate w_x which contributes to airborne volumetric activity through coolant has been simulated. The volumetric activity of ^{131}I inside the containment for various values of mixing rate w_x from 0.005 to 1.0 s^{-1} has been numerically simulated. It has been observed that a slight change in the value of w_x results in astonishing variation in airborne activity. The results are shown in **Figure 7**. For $w_x = 1.0 \text{ s}^{-1}$, highest magnitude has been observed with a shift of peak toward lower timescale value. However, the ^{131}I has been observed to reduce almost linearly with mixing rate 1.0 s^{-1} (**Figure 7**). The in-containment volumetric concentration of ^{131}I increases and then starts decreasing gradually for mixing rate value from $w_x = 1.0$ to 0.03 s^{-1} ; however, it remains constant for mixing rate 0.005 s^{-1} . This is because of trivial mixing of iodine in the coolant. A higher magnitude of ^{131}I activity has been observed with higher mixing rate.

5.4. Fission product activity with spray system

The primary purpose of the spray system is to mitigate the FP exposure to the environment and to maintain the containment integrity. In this work, we have studied the effect of the spray system in mitigating the radioactive masses (gaseous and particles) released during in-vessel release phase. During loss of coolant accident, the temperature and pressure inside the containment start raising. It reached to 80° pressure and reached to 7.533 psi within few minutes. The simulation has been carried out by assuming the containment temperature at 80°C and pressure at 7.533 psi and the spray with pH 5.0 and 9.5 and the alkaline spray. The spray system is found

to have minimum effect on noble gases and reduces iodine and other radioactive particles effectively. The spray system is started at 100, 500, 1000, and 1500 s after the release time. The effect of spray system activation time on noble gases and iodine is shown in **Figure 8**.

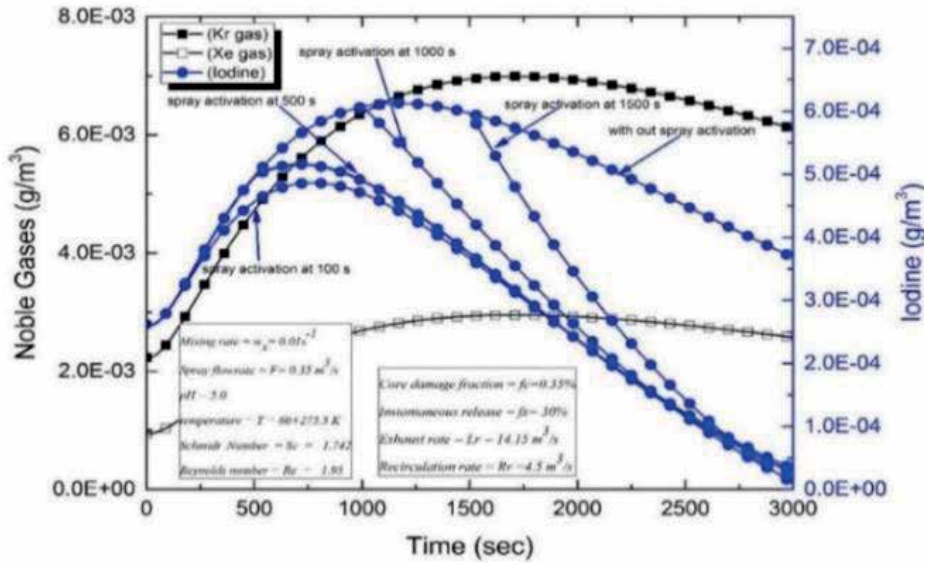


Figure 8. Radioactive noble gases and iodine release during in-vessel release phase with mixing rate $w_x = 0.01 \text{ s}^{-1}$.

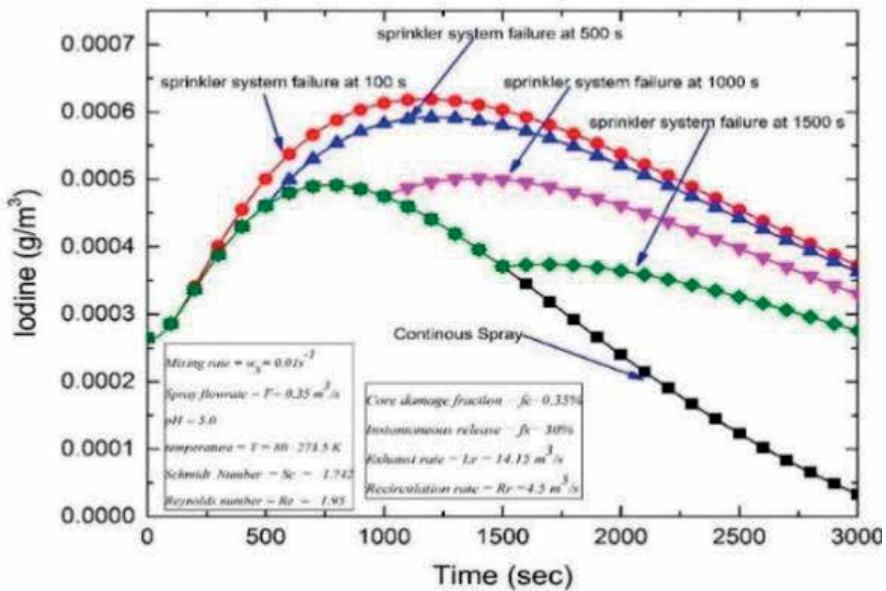


Figure 9. Iodine with containment spray system failures during in-vessel release phase with mixing rate $w_x = 0.01 \text{ s}^{-1}$.

The iodine concentration first increases exponentially in containment and immediately starts reducing with the activation of the spray system. This is because of competition between the continuous source of radioactive iodine (P_i) coming from reactor pressure vessel (RPV) and removal of iodine with the containment spray system (Figure 8). The volumetric iodine starts reducing exponentially after the activation of the spray system. The droplet size has been assumed 800 microns for the simulation. However, the noble gases are observed to be unaffected with a spray system. The response of airborne iodine to the failure of the spray system is depicted in Figure 9.

Figure 9 indicates that the premature failure of the system ($t = 100$ s) does not affect the airborne iodine concentration. The slight decrease in airborne iodine concentration is seen if the spray system is failed or malfunctions at 500 s. However, the spiracles are operated for a longer period, for example, 1000–1500 s during the in-vessel phase. The airborne concentration of iodine is reduced significantly. The failure of the system at 1000 and 15,000 s caused the regaining of airborne iodine (Figure 9).

5.5. Droplet diameter and pH effect on in-containment FPA

The droplet collection efficiency of spray also depends on the containment atmospheric temperature, pressure and spiracle pH value. The effect of spray water pH value and an alkaline spray has been simulated. The results are depicted in Figure 10. The results showed that the higher pH spray solution (pH 9.5) and alkaline solution ($\text{Na}_2\text{S}_2\text{O}_3$) have similar removal characteristics for airborne iodine. The iodine removal rates of Boric (pH 5.0), NaOH (pH 9.5) and alkaline solution with different droplet sizes are shown in Figure 11.

The removal rate has been seen to decrease exponentially with the increase in droplet size for an alkaline solution (Figure 11). However, for a spray solution with pH value 5.0, the removal rate decreases in a linear manner. The in-containment volumetric mass under the atmospheric

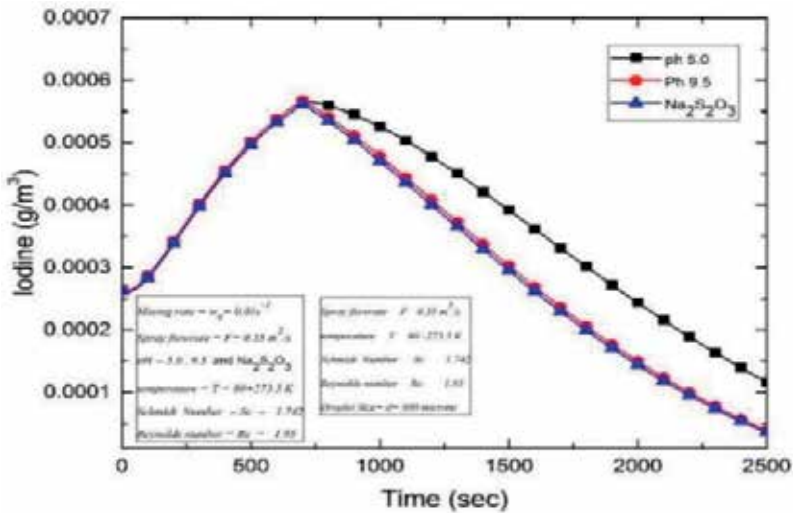


Figure 10. Response of radioactive iodine for spray pH value during in-vessel release phase with mixing rate $w_x = 0.01 \text{ s}^{-1}$.

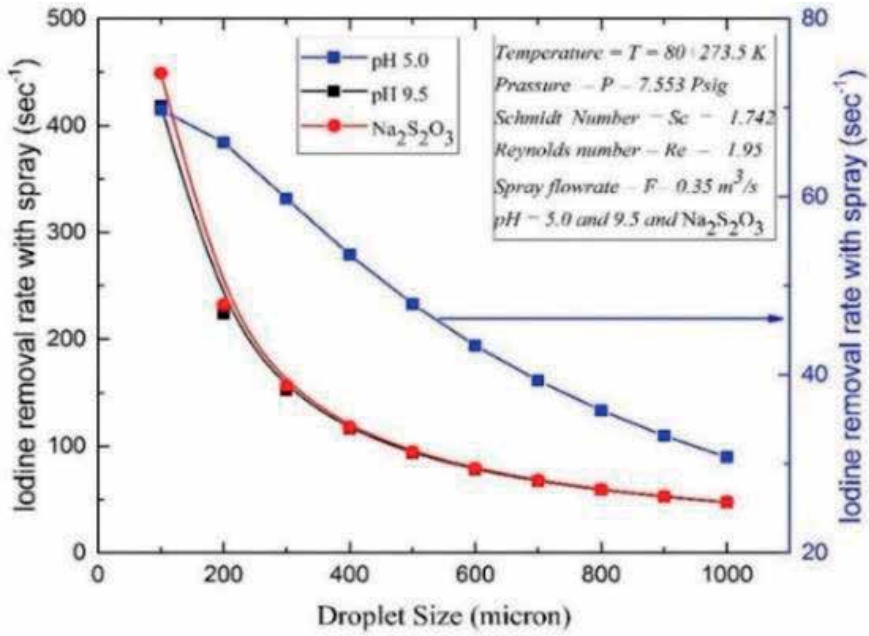


Figure 11. Droplet removal rate for iodine at pH 5.0, 9.5 and with alkaline spray solution ($\text{Na}_2\text{S}_2\text{O}_3$).

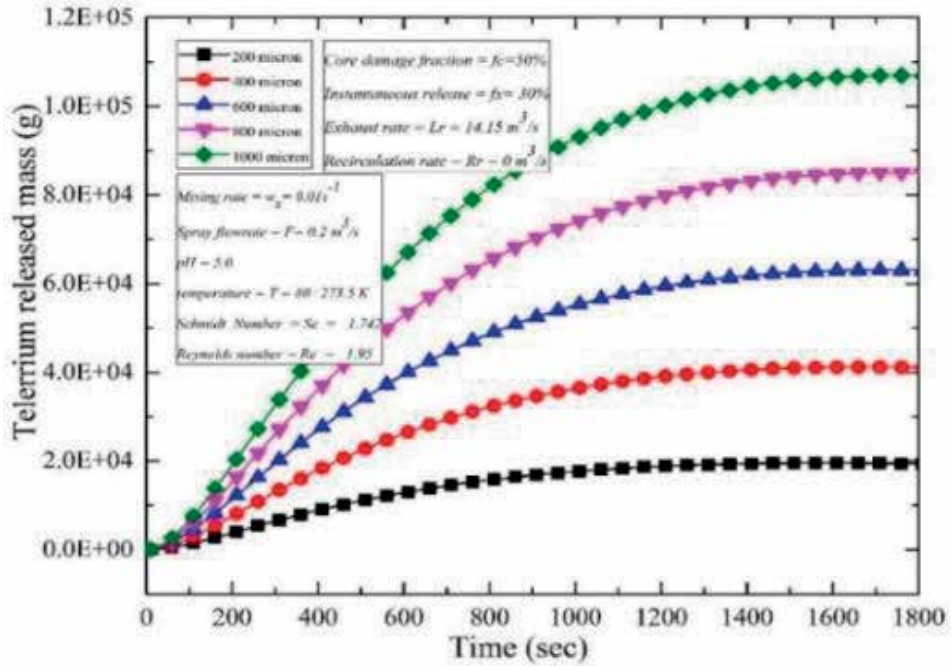


Figure 12. Tellurium inventory during in-vessel release phase with mixing rate $w_x = 0.01 \text{ s}^{-1}$.

conditions with a temperature of 80°C and 7.533 Psi. Assuming 35% core damage and 20% burst release. The rest of mass is assumed to release along with the coolant with mixing rate $w_x = 0.01\text{ s}^{-1}$. The mass concentration of tellurium is simulated for droplet sizes (100–1000 microns). The containment spray system is activated with the initiation of an accident with a constant flow rate of $0.2\text{ m}^3/\text{s}$.

Simulation results showed that the droplet size is quite effective to reduce the airborne FPs. It has been observed that the concentration of airborne tellurium decreases with a decrease in droplet size (**Figure 12**). The peak concentration in Tellurium mass reaches to a maximum concentration at a longer time with the higher droplet diameter. The magnitude of maximum concentration has been found the approximately inverse square of droplet diameter ($1/d^2$). The containment spray system removal rate for iodine versus droplet diameter is depicted in **Figure 11**. The maximum removal rate has been found 452 s^{-1} with alkaline solution spray with a droplet size of 100 micrometers. The removal rate is found to decrease exponentially as the droplet diameter increases. 44.7 s^{-1} removal rate has been seen for 1000-micron diameter droplet size for pH 9.5 and alkaline spray solution with spray flow rate $0.35\text{ m}^3/\text{s}$. The gas phase and liquid phase coefficients play a vital role in absorption efficiency. Both gas- and liquid-phase mass transfer coefficients (K_G and K_L) decrease drastically with an increase in droplet size (**Figure 13**). However, the gas- and liquid-phase mass transfer coefficients (K_G and K_L) are also related to the inverse square of droplet diameter ($1/d^2$).

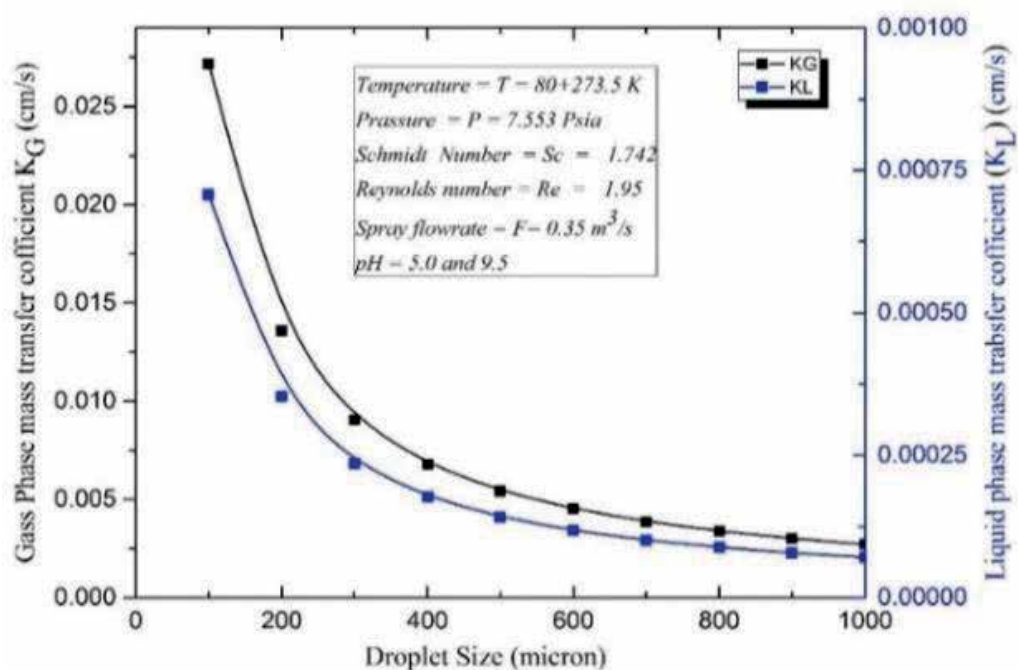


Figure 13. A comparison of gas-phase mass transfer coefficient (K_G) and liquid-phase mass transfer coefficients (K_L) for elemental iodine constant pH = 5.0, $T = 80^{\circ}\text{C}$.

6. Conclusion

This chapter has presented the numerical simulation of FP activity inside the reactor containment building under LOCA. The numerical simulation of in-containment FPs against the mixing rate, puff release, droplet diameter, spray pH value and spray performance has been simulated. The results indicate that the mixing rate of FPs in coolant significantly affects airborne FP activity inside the containment. The higher pH spray solution (9.5) and spray with sodium thiosulfate ($\text{Na}_2\text{S}_2\text{O}_3$) have observed similar scrubbing properties. The droplet size is significantly important for removal of FP. There is a higher tendency of FP to interact with airborne particles (**Figure 11**) with, due to their higher values of liquid- and gas-phase mass transfer coefficients (K_L and K_G) (**Figure 13**). Therefore, the acceptance criteria of droplet size have been suggested between 600 and 800 microns, with pH value higher than 7.0 which delivers higher removal rate. The earlier the containment spray system has operated, the airborne concentration will be minimum (**Figure 8**). However, the delay operation caused the higher airborne concentration of radioactive mass. It has also been observed if the containment spray system is failed during in-vessel phase a regain will be caused in radioactive mass (**Figure 9**). Based on our work, we are suggesting 600–8500 microns mean droplet diameter containment spray system should be used to get maximum radiation hazard safety. Moreover, from our results, we can conclude that the spray system should be operated within 500 s after the accident and should be operated more than 3000 s (whole in vessel phase). The uncertainties in simulated results depend on generally available data in the literature.

Acknowledgements

This article was funded by the deanship of scientific research (DSR) at King Abdul Aziz University, Jeddah. The author, therefore, acknowledges with thanks DSR for technical and financial support.

Author details

Khurram Mehboob* and Mohammad Subian Aljohani

*Address all correspondence to: khurramhrbeu@gmail.com

Department of Nuclear Engineering, Faculty of Engineering, King Abdul Aziz University, (KAU), Jeddah, Saudi Arabia

References

- [1] Rahim FC, Rahgoshay M, Mousavian SK. A study of large break LOCA in the AP1000 reactor containment. *Progress in Nuclear Energy*. 2012;**54**:132-137

- [2] Huang GF, et al. Study on mitigation of in-vessel release of fission products in severe accidents of PWR. *Nuclear Engineering and Design*. 2010;**240**(11):3888-3897
- [3] Denning RS, et al. Radionuclide Release Calculations for Selected Severe Accident Scenarios. NRC NUREG/CR-4624, 1986
- [4] Lewis BJ, et al. Fission product release mechanisms during reactor accident conditions. *Journal of Nuclear Materials*. 1999;**270**:21-38
- [5] Iglesias FC, et al. Fission product release mechanisms during reactor accident conditions. *Journal of Nuclear Materials*. 1999;**270**:21-38
- [6] Mehboob K, et al. Source term evaluation of two loop PWR under hypothetical severe accidents. *Annals of Nuclear Engineering*. 2012;**50**:271-284
- [7] Tigeras A, Bachet M, Catalette H, Simoni E. PWR iodine speciation and behavior under normal primary coolant conditions: An analysis of thermodynamic calculations, sensibility evaluations and NPP feedback. *Progress in Nuclear Energy*. 2011;**53**:504-515
- [8] Soffer L, et al. Accident Source Terms for Light-Water Nuclear Power Plants. NUREG-1465, U.S. Nuclear Regulatory Commission 1995
- [9] Dehjourian M, et al. Effect of spray system on fission product distribution in containment during a severe accident in a two-loop pressurized water reactor. *Nuclear Engineering and Technology*. 2016;**48**:975-981
- [10] Rosen M, Jankowski M. Reassessing releases: A closer look at source term. IAEA Bulletin. Special Report. 1985:43-46
- [11] Ducros G, et al. Fission product release under severe accidental conditions general presentation of the program and synthesis of VERCORS 1-6 results. *Nuclear Engineering and Design*. 2001;**20**(2):191-203
- [12] Lewis BJ, et al. Overview of experimental programs on core melt progression and fission product release behaviour. *Journal of Nuclear Materials*. 2008;**380**(1-3):126-143
- [13] Luis E, Herranz B. Clement in-containment source term key insights gained from a comparison between the PHEBUS-FP programme and the US-NRC NUREG-1465 revised source term. *Progress in Nuclear Energy*. 2010;**52**(5):481-486
- [14] Giraulta N, et al. LWR severe accident simulation: Iodine behavior in FPT2 experiment and advances on containment iodine chemistry. *Nuclear Engineering and Design*. 2012;**243**:371-392
- [15] Hastet T, et al. Phébus FPT3: Overview of main results concerning the behaviour of fission products and structural materials in the containment. *Nuclear Engineering and Design*. 2013;**261**:333-345
- [16] Mäkynen JM, et al. Experimental studies on hygroscopic aerosol behaviour in LWR containment conditions. *Journal of Aerosol Science*. 1994;**25**(Suppl. 1):247-248

- [17] Mäkynen JM, et al. AHMED experiments on hygroscopic and inert aerosol behaviour in LWR containment conditions: Experimental results. *Nuclear Engineering and Design*. 1997;**178**(1-2):45-59
- [18] Schwarz M, et al. PHEBUS FP: A severe accident research programme for current and advanced light water reactors. *Nuclear Engineering and Design*. 1999;**187**:47-69
- [19] Girault N, et al. Towards a better understanding of iodine chemistry in RCS of nuclear reactors. The 2nd European Review Meeting on Severe Accident Research (ERMSAR-2007) Forschungszentrum Karlsruhe GmbH (FZK), Germany, 12–14 June 2007
- [20] Di Giuli M, et al. SARNET benchmark on Phébus FPT3 integral experiment on core degradation and fission product behaviour. *Annals of Nuclear Energy*. 2016;**93**:65-82
- [21] Kljenak, et al. Thermal-hydraulic and aerosol containment phenomena modelling in ASTEC severe accident computer code. *Nuclear Engineering and Design*. 2010;**240**(3):656-667
- [22] Haste T, et al. MELCOR/MACCS simulation of the TMI-2 severe accident and initial recovery phases, off-site fission product release and consequences. *Nuclear Engineering and Design*. 2006;**236**(10):1099-1112
- [23] Slaga AC, et al. MAAP4.0.7 severe accident source term analysis, In Proceedings of the International Congress on Advances in Nuclear Power Plants (ICAPP ' 08), 2008. Anaheim, Calif, USA, Paper 8201
- [24] Singh M, et al. A numerical methodology for estimation of volatile fission products release from nuclear fuel. *Nuclear Engineering and Design*. 2017;**(323)**:338-344
- [25] Lewis BJ, et al. Fission product release modelling for application of fuel-failure monitoring and detection - an overview. *Journal of Nuclear Materials*. 2017;**489**:64-83
- [26] Lewis BJ. A generalized model for fission product transport in the fuel-to-sheath gap of defective fuel elements. *Journal of Nuclear Materials*. 1990;**175**(3):218-226
- [27] Koo YH, Sohn DS, Yoon YK. Release of unstable fission products from defective fuel rods to the coolant of a PWR. *Journal of Nuclear Materials*. 1994;**209**(3):248-258
- [28] Avanov AS. The model of the fission gas release out of porous fuel. *Ann.Nucl. Eney*. 1998;**25**(15):1275-1280
- [29] Tucker MO, White RJ. The release of fission products from UO₂ during irradiation. *Journal of Nuclear Materials*. 1979;**87**(1):1-10
- [30] Awan SE, et al. Sensitivity analysis of fission product activity in primary coolant of typical PWRs. *Progress in Nuclear Energy*. 2011;**53**(3):245-249
- [31] Javed Iqbal M, et al. Kinetic simulation of fission product activity in primary coolant of typical PWRs under power perturbations. *Nuclear Engineering and Design*. 2007;**237**(2):199-205
- [32] USNRC. Alternative Radiological Source Terms for Evaluating Design Basis Accidents at Nuclear Power Reactors, Regulatory Guide 12000. p. 183

- [33] USNRC, Reactor Safety Study. An Assessment of Accident Risk in U.S. Commercial Nuclear Power Plants, WASH-1400, United States Nuclear Regulatory Commission 1975
- [34] Henry RE, TMI-2: A Text Book in Severe Accident Management, MISD Professional Development Workshop, ANS/ENS international meeting, 2007
- [35] Jak K. Nuclear Power Plant Modelng and Steam Generator Stability Analysis. PhD Thsis. The University of Michigan; 1981
- [36] Mehboob K, et al. Numerical simulation of radioisotope's dependency on containment performance for large dry PWR containment under severe accidents. Nuclear Engineering and Design. 2013;**262**:435-445
- [37] El-Jaby, et al. A general model for predicting coolant activity behavior for fuel failure monitoring analysis. Journal of Nuclear Materials. 2010;**399**:87-100
- [38] Mehboob K, Aljohani SM. Modeling and simulation of radio-iodine released inside the containment as result of an accident. Progress in Nuclear Energy. 2016;**88**:75-87

Study of the Numerical Diffusion in Computational Calculations

Despoina P. Karadimou and
Nikos-Christos Markatos

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.75660>

Abstract

The numerical simulation of fluid flow and heat/mass-transfer phenomena requires the numerical solution of the Navier-Stokes and energy-conservation equations coupled with the continuity equation. Numerical or false diffusion is the phenomenon of inserting errors in the calculations that compromise the accuracy of the computational solution. The Taylor series analysis that reveals the truncation/discretization errors of the differential equations terms should not be termed as false diffusion. Numerical diffusion appears in multi-dimensional flows when the differencing scheme fails to account for the true direction of the flow. Numerical errors associated with false diffusion are investigated via two- and three-dimensional problems. A numerical scheme must satisfy necessary criteria for the successful solution of the convection-diffusion formulations. The common practice of approximating the diffusion terms via the central-difference approximation is satisfactory. Attention is directed to the convection terms since these approximations induce false diffusion. The equations of all the conservation equations in this study are discretized by the finite volume method.

Keywords: false diffusion, numerical dispersion, computational errors, finite volume method

1. Introduction

Numerical diffusion is a significant source of error in numerical solution of conservation equations and can be separated into two components namely cross-stream and streamwise numerical diffusion. The former occurs when gradients in a convected quantity exist perpendicular to the flow and the direction of the flow is oblique to the grid lines, i.e. due to

the multi-dimensional nature of the flow. The latter happens when gradients in a convected quantity exist parallel to the flow even in one-dimensional situations [1–3].

According to Vahl Davis et al. [4] a theoretical approximation for the false-diffusion term in two-dimensional geometries is:

$$\Gamma_{false} = \frac{\rho U \Delta x \Delta y \sin 2\theta}{4(\Delta y \sin^3 \theta + \Delta x \cos^3 \theta)} \quad (1)$$

where U is the resultant velocity and θ is the angle, which lies between 0 and 90°, made by the velocity vector with the x -direction. Expression (1) reveals that there is no false-diffusion present for $\theta = 0$ and 90°, and that it is at its maximum at $\theta = 45^\circ$.

The first-order accurate upwind-difference scheme is more stable than the second-order accurate central-difference scheme or higher-order accurate schemes, e.g., QUICK at high grid Peclet numbers. Non-linear schemes, e.g., van LEER scheme appear highly stable but they do not deal with the flow direction inclination at all [5–8].

One way to overcome diffusion errors is to use an upwind approximation which essentially follows the streamlines. This approach, originally derived by Raithby [9], is formally called the skew-upwind differencing scheme. The CUPID (Corner UPwInDing) scheme retains the general objectives of the Raithby approach but uses an entirely different 2D formulation that eliminates the shortcomings of the original scheme. The SUCCA (Skew Upwind Corner Convection Algorithm) scheme is based on a simplified formulation of the CUPID scheme [10, 11].

In this chapter the performance of various numerical schemes is compared as far as the accuracy of the numerical flow prediction is concerned. Various cases of numerical problems that include one- and two-phase flows, heat and mass transfer, geometry of two and three-dimensions, grid of Cartesian and curvilinear coordinate system, straight and inclined inflows are investigated and the conclusions are presented. Among the numerical schemes that are evaluated is the SUPER (Skew Upwind and cornER algorithm) numerical scheme that takes into account the 3D flow orientation and follows all the rules for the successful solution of the equations [12].

Numerical dispersion comes about when the convective scheme used is unstable and large gradients are present [3]. The use of the upwind or hybrid numerical scheme ensures the stability of the calculations but the first-order accuracy makes them prone to streamwise numerical diffusion errors. Higher-order schemes involve more neighbour points and reduce the streamwise false-diffusion by bringing in a wider influence. These schemes are based on one-dimension formulations that do not take into account the flow direction which makes them unstable at high grid Peclet numbers [13]. Furthermore numerical schemes that take into account the flow direction become unstable at high grid Peclet numbers in regions of high convection discontinuities. The phenomenon of the numerical dispersion causes spatial oscillations that lead to many numerical errors and non-physical results [1, 2].

Generally the numerical schemes should follow some basic rules in order to be able to calculate a successful solution:

- a. The formulation of the scheme should ensure the flux consistency at the cell faces in adjacent control volumes satisfying the conservative property.
- b. Conservation principle should be retained for all the variables and in the whole domain.
- c. The scheme should recognize the relationship between the magnitude of the Peclet number and the directionality of influencing known as transportiveness property.
- d. The influence coefficients in the finite volume equations should be always positive.
- e. The neighbouring coefficients for the solution control volume of central node P should obey the following relation: $\sum a_{nb} = a_p$, so as to be consistent with the differential transport equation.
- f. Satisfying the Scarborough criterion $\frac{\sum |\alpha_{nb}| \leq 1 \cdot \text{at} \cdot \text{all} \cdot \text{nodes}}{|a_p| < 1 \cdot \text{at} \cdot \text{least} \cdot \text{one} \cdot \text{node}}$ ensures that the resulting matrix of coefficients is diagonally dominant which is a sufficient condition for a converged iterative method.
- g. Linearizing the source term with a negative slope, so that the incidence of unbounded solutions is reduced.
- h. The formulation of the scheme should be easy to implement without increasing the computational cost.

The method employed in this study follows the so called “control volume” approach which is developed by formulating the governing equations on the basis that mass, heat and momentum fluxes are balanced over control volumes. If there are two phases the control volume can be regarded as containing a volume fraction of each phase (r_i) that obeys the following relation: $r_1 + r_2 = 1$. Each phase is treated as a continuum in the control volume under consideration. The phases share the control volume and they may, as they move within it, interpenetrate [14, 15].

1.1. Mathematical modeling

The general form of the differential transport equations in the Cartesian coordinate system is described as follows [14, 15]:

$$\frac{\partial}{\partial t}(r_i \rho_i \phi_i) + \text{div}(r_i \rho_i \vec{u}_i \phi_i - r_i \Gamma_{\phi,i} \text{grad} \phi_i) = S_{\phi,i} \quad (2)$$

where ϕ_i the dependent variable of each phase i ; r_i the volume fraction of each phase, (m^3/m^3); ρ_i the density of each phase, (kg/m^3); u_i the velocity vector of each phase, (m/s); $\Gamma_{\phi,i}$ within-phase diffusion coefficient; and $S_{\phi,i}$ within-phase volumetric sources, ($\text{kg}/\text{m}^3 \text{ s}$).

The dependent variables solved for are:

- a. Pressure that is shared by the two phases, P (Pa)
- b. Volume fractions of each phase, r_i (m^3/m^3)
- c. Three components of velocity for each phase, u_i, v_i, w_i (m/s)
- d. Turbulence kinetic energy and dissipation rate of turbulence for the first phase, k and ε (m^2/s^2)
- e. Temperature (K), enthalpy (J/kg)

The phase volume fraction equation is obtained from the continuity equation:

$$\frac{\partial(r_i \rho_i)}{\partial t} + \text{div}(r_i \rho_i \vec{u}_i) = S_{\phi,i} \quad (3)$$

where r_i = phase volume fraction, (m^3/m^3); ρ_i = phase density, (kg/m^3); \vec{u}_i = phase velocity vector, (m/s); and $S_{\phi,i}$ = net rate of mass entering phase i from phase j , ($\text{kg}/\text{m}^3 \text{ s}$), if there is phase change.

The discretized by the finite volume method form of all the conservation equations is solved by the SIMPLEST and IPSA algorithms embodied in the Computational Fluid Dynamics (CFD) code PHOENICS [16].

2. Numerical problems

2.1. Inclined flow of a scalar quantity

The physical problem presented in this section describes the transport of a scalar quantity (C1) in a two-dimensional geometry. The flow of C1 is inclined to the grid lines ($\theta = 45^\circ$) and the natural diffusion is assumed equal to zero. This physical problem is considered as a benchmark for comparing the performance of various numerical schemes.

The boundary conditions applied to the grid are the following:

Air enters the domain from the west side with stable mass flow rate ($\rho \cdot u$) $\text{kg}/\text{m}^2 \text{ s}$, uniform values of the two components of velocity ($u_1 = v_1 = 1 \text{ m/s}$) and uniform value of the convected quantity (C1) ($1 \cdot \text{kg}/\text{m}^3$).

Air enters the domain also from the south side with stable mass flow rate ($\rho \cdot u$) $\text{kg}/\text{m}^2 \text{ s}$, uniform values of the two components of velocity ($u_1 = v_1 = 1 \text{ m/s}$) and uniform value of the convected quantity (C1) ($1 \cdot \text{kg}/\text{m}^3$).

At the two outlets (north and east side) of the domain the air is supposed to exhaust at an environment of fixed uniform pressure.

The dimensions of the domain are $1 \times 1 \text{ m}$. The numerical results presented are based on the independent grid of 33×33 cells. For the discretization of the convected term in the transport equation of the scalar quantity C1 the numerical schemes: (a) UPWIND, (b) van LEER, (c) SUCCA, (d) SUPER are applied [6, 8, 11, 12].

In **Figure 1** the velocity vector distribution predicted by the SUPER numerical scheme is presented.

In **Figure 2** the vertical (concentration) distribution of the scalar quantity C1 in the middle of the domain is presented.

The vertical distribution of the scalar quantity C1 that is transferred by air with inclined direction is predicted more abruptly by the numerical schemes (SUCCA, SUPER) that take into account the

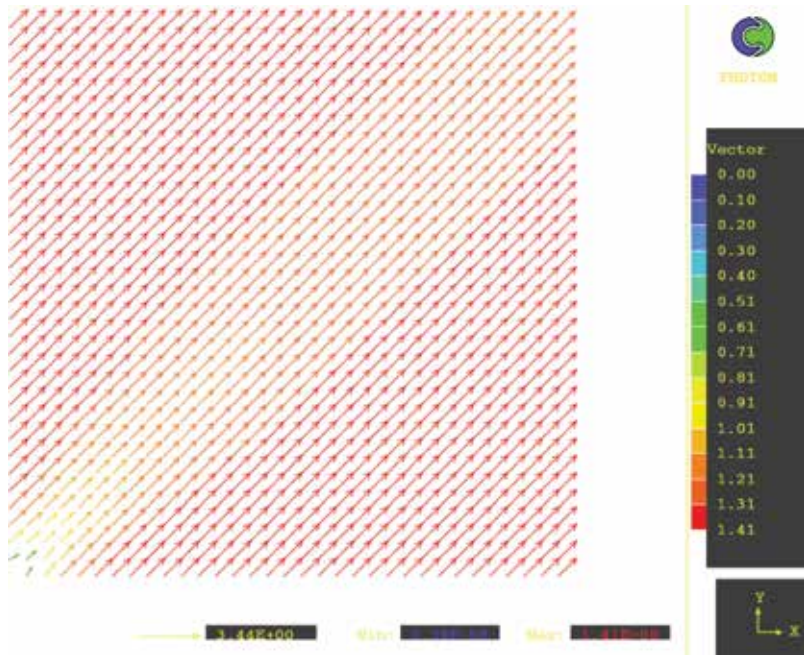


Figure 1. Air velocity vector distribution predicted by the SUPER scheme.

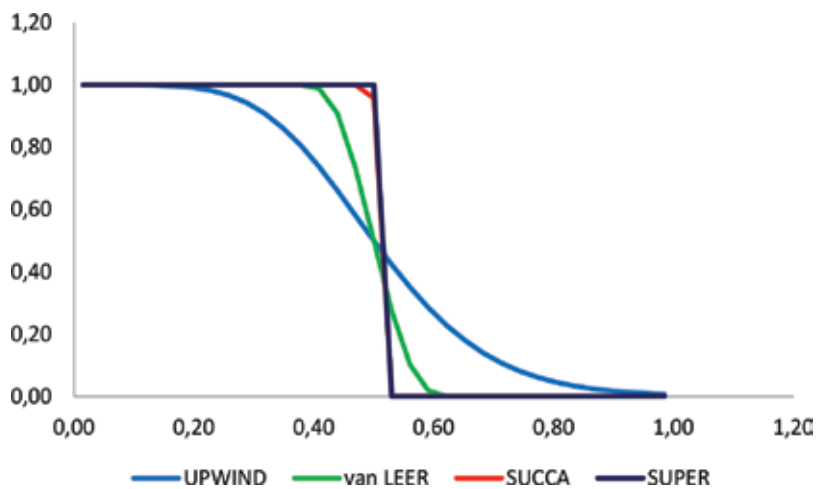


Figure 2. Vertical distribution of the scalar quantity C1 in the middle of the domain applying the numerical schemes: (a) UPWIND, (b) van LEER, (c) SUCCA, (d) SUPER.

flow direction than the conventional numerical scheme (UPWIND) and the non-linear numerical scheme (van LEER). It is concluded that the numerical diffusion errors are minimized when the discretization schemes (SUCCA, SUPER) are applied to predict the transport of the scalar quantity C1 when the air-flow direction appears the major inclination angle 45° to the grid lines. The performance of the van LEER numerical scheme is also satisfactory in predicting the scalar concentration distribution. The conventional UPWIND numerical scheme that does not take into account the phenomenon of false-diffusion presents the poorest accuracy.

2.2. Tubular flow of a scalar quantity and heat conduction in the radial direction

A physical problem used in this study for assessing the numerical schemes is the numerical prediction of the energy transfer in a triple tube heat exchanger [17]. The heat exchanger has a three-dimensional curvilinear geometry similar to a circular cylinder. Three fluids are being considered to flow inside the domain that is separated by solid. Chilled water (10°C) flows in the inner tube, hot water (70°C) flows in the inner annulus and normal tap water (18°C) flows in the outer annulus. The three fluids follow a co-current flow. The numerical model is first validated and then used to compare the performance of various numerical schemes.

The boundary conditions applied for solving the above physical problem are the following:

Water with uniform properties (temperature, density, thermal capacity, kinematic viscosity) enters the three different inlets of the heat exchanger. The heat exchanger is separated to three components that are filled with water of different properties and the solid parts that prevent the water vertical flow. Heat is exchanged between the fluids through conduction. The solid is steel at 27°C of uniform properties (density, viscosity, specific heat, thermal conductivity, thermal expansion coefficient, compressibility). At the three outlets of the heat exchanger the boundary condition of zero mass flow is applied.

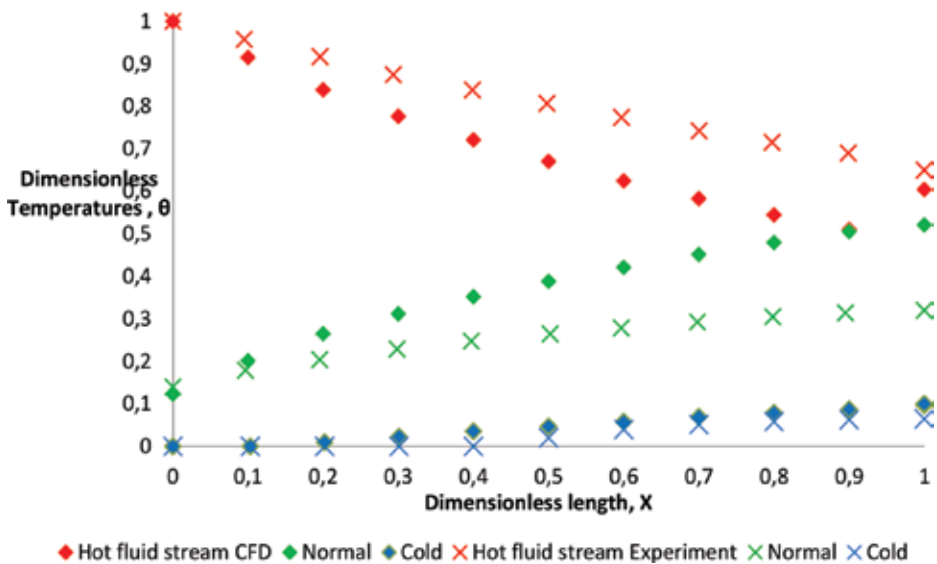


Figure 3. Local temperature variation of the co-current three fluid streams along the length of the heat exchanger.

In **Figure 3** the numerical prediction of the temperature distribution along the heat exchanger is validated against experimental data. The average% relative error of the numerical prediction does not exceed the value of 20%.

In **Figure 4** the numerical results of the enthalpy distribution at the lateral plane in the middle of the heat exchanger applying three different numerical schemes (HYBRID, van LEER, SUPER) [6, 8, 12] is presented.

Comparing the performance of the numerical schemes in predicting the transfer of the scalar quantity H1 it is concluded that the numerical prediction is equally satisfactory for the three discretization schemes.

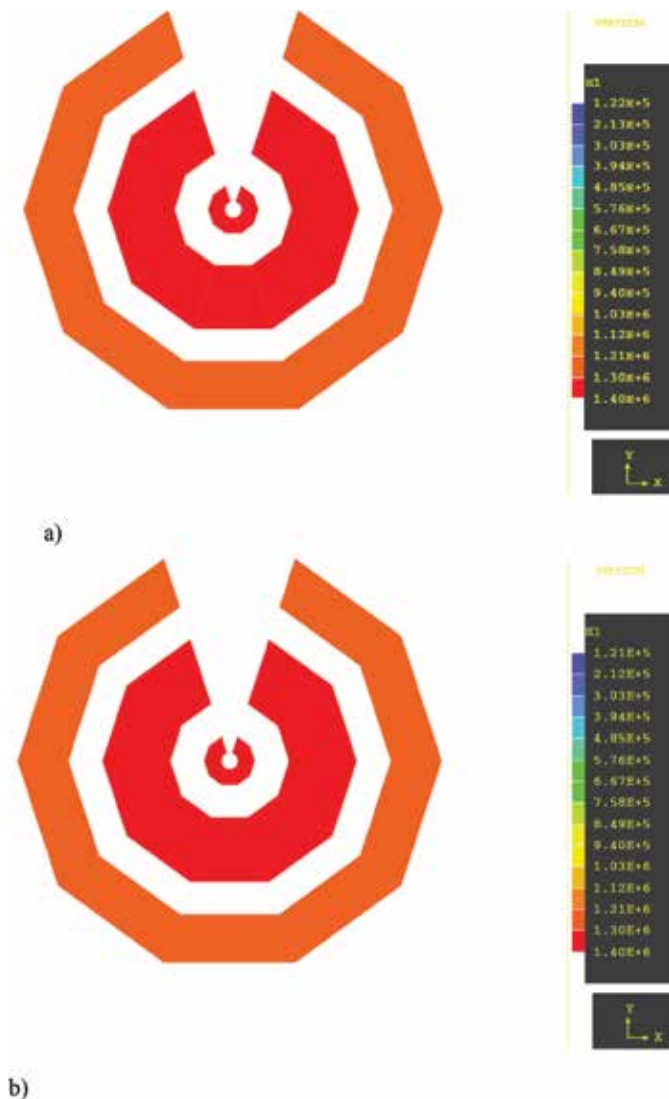


Figure 4. Enthalpy distribution at the lateral plane of the heat exchanger applying: (a) the HYBRID, (b) the van Leer, (c) the SUPER numerical scheme.

2.3. Airflow in a backward facing step

A physical problem appropriate for assessing the accuracy of various numerical schemes is the prediction of the airflow velocity vectors distribution in a backward facing step.

The separation of the airflow, the recirculation zone formed forward the step and the reattachment length are the main characteristics of the flow in this geometry. The numerical solution of the airflow field in a backward facing step is of major concern due to the difficulty in accurate prediction of these complex phenomena. The flow field may become more complex in the presence of particles.

In this study the two-dimensional backward facing step of Benavides and Wachem [18] is investigated.

In **Figures 5** and **6** the velocity flow field is presented applying various numerical schemes (UPWIND, HYBRID, van LEER, SUPER) [6, 8, 12] and validated against experimental data [18].

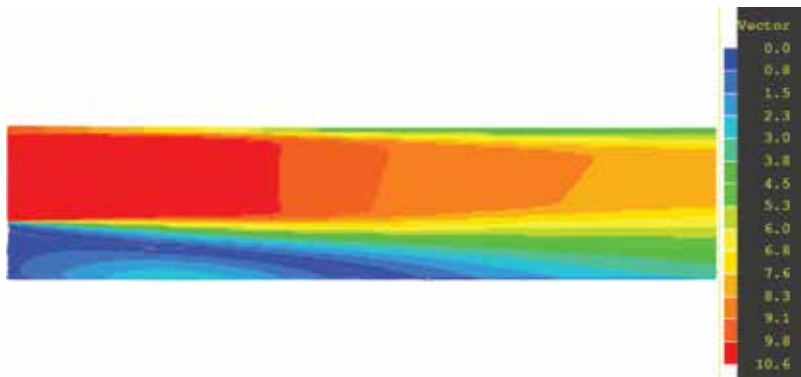


Figure 5. Velocity flow field applying various numerical schemes (upwind, hybrid, van Leer, super).

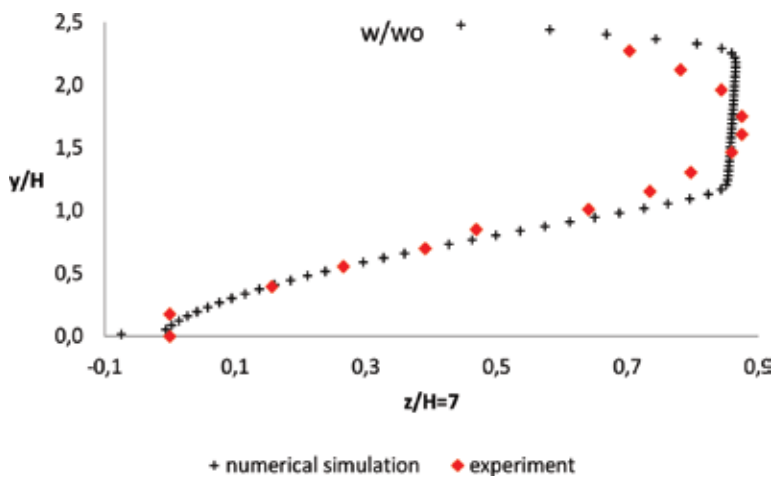


Figure 6. Vertical velocity distribution applying various numerical schemes (UPWIND, HYBRID, van LEER, SUPER) and comparison with the experimental data.

According to the numerical results that agree well with the experimental data the performance of the numerical schemes (UPWIND, HYBRID, van LEER, SUPER) is equally satisfactory. The same conclusion is derived for the case of air-particles flow in a backward facing step.

2.4. Inclined air-particles flow

A physical problem that has been used for assessing the performance of various numerical schemes is the dispersed air-particles flow in a model room geometry (**Figure 7**) [19]. A two-phase Euler-Euler mathematical model is applied to calculate the oblique inflow in the interior of the internal space. The accuracy of the four numerical schemes: (a) the conventional first-order UPWIND numerical scheme, (b) the second-order HYBRID scheme, (c) the non-linear van LEER and (d) the flow-oriented SUPER scheme [6, 8, 12] are compared in the case of inclined inflow ($\theta = 45^\circ$).

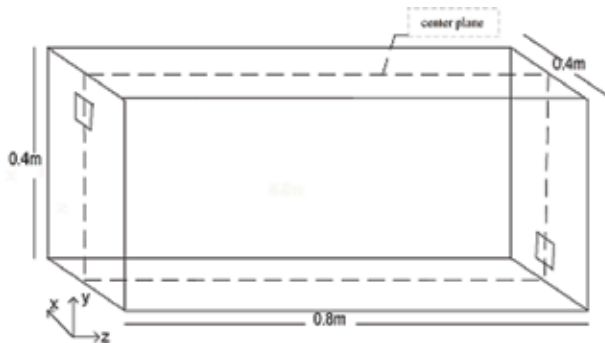


Figure 7. Geometry of the model room.

2.4.1. Boundary conditions

At the inlet the mass flow rate of each phase is multiplied by its volume fraction. The dispersed turbulent air-particles flow enters the room with uniform velocity (0.225 m/s) in an inclined direction. Turbulence is modeled by the RNG k- ϵ model [20]. Particles are assumed to be transported and dispersed due to turbulence of the carrier fluid (air).

The turbulence kinetic energy of the air-phase that is applied at the inlet is defined as [13, 21] $k_{in} = \frac{3}{2} (U_{avg} T_i)^2$ where U_{avg} is the mean air inlet velocity and T_i the turbulence intensity, considered as 6%. The dissipation rate is given by $\epsilon = C_\mu^{3/4} k^{3/2} / \ell$, where ℓ the turbulence length scale is taken as $\ell = 0.07 d_h$ (d_h hydraulic diameter of the duct) and $C_\mu = 0.0845$ an empirical constant of the turbulence model. The mathematical model uses the logarithmic "wall functions" near the solid surfaces.

At the outlet both air phases are supposed to exhaust at an environment of fixed uniform pressure. At the walls the no-slip and no-penetration condition is applied for both phases.

2.4.2. Results

Figure 8a–c present the vertical w_1 velocity distribution at the longitudinal plane of the domain at distance 0.2, 0.4, 0.6 m from the supply inlet.

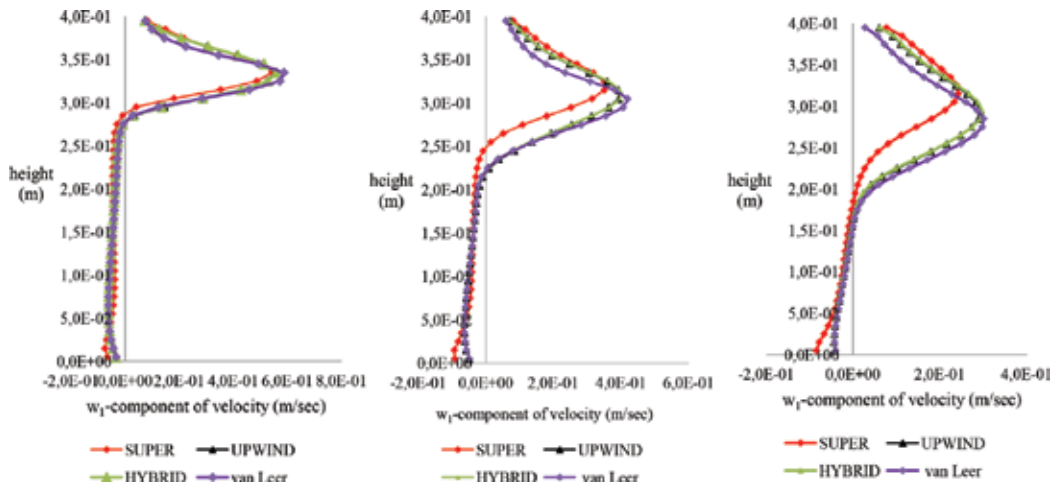


Figure 8. Vertical w_1 velocity distribution at the longitudinal plane of the domain and (a) 0.2 m from the inlet; (b) 0.4 m from the inlet; and (c) 0.6 m from the inlet

As far as the accuracy of the four numerical schemes in the case of inclined inflow ($\theta = 45^\circ$) is concerned the higher-order and non-linear schemes (hybrid and van Leer numerical scheme) present a similar performance with the first-order upwind numerical scheme and a different performance than the flow-oriented scheme (SUPER scheme). The vertical w_1 velocity distribution predicted by the SUPER scheme presents a more abrupt and accurate profile due to the successful minimizing of the false-diffusion errors.

2.5. Heat and mass transfer

A physical problem that has been used to evaluate the performance of the numerical accuracy is the water-vapor condensation of humid air in the three dimensional geometry of a real-scale indoor space. A two-phase flow Euler-Euler mathematical model has been developed, wherein the humid air and water droplets are being treated as separate phases. The two phases exchange momentum and energy and, as the temperature drops below the dew point of humid air, mass transfer and phase change of water vapor to liquid takes place. The flow of humid air inside the room is buoyancy driven in the temperature range of 290–303 K. The properties of humid air (enthalpy, relative humidity, concentration of water vapor, saturation vapor pressure) vary with the temperature [22]. The dimensions of the domain are: width (X) \times height (Y) \times length (Z) = 4.0 \times 4.0 \times 8.0 m.

2.5.1. Boundary and initially conditions

The interior of the domain is filled with humid air of 303 K and no liquid phase. The temperature on the surface of the walls is 303 K and on the surface of the cold bottom is 290 K. The density of humid air is 1.16 kg/m³ and the water droplets have a constant density (996 kg/m³) and specific heat (4190.0 J/kg \cdot K) at the dew point temperature. The initial pressure inside the

domain is equal to the atmospheric pressure. The initial relative humidity condition tested is 90%. The overall water vapor mass of humid air at the initial temperature 303 K is taken from the psychrometric chart [23].

2.5.2. Results

In **Figure 9** the vertical temperature distribution in the middle of the domain at the height (0–4 m) is presented.

The temperature of humid air near the floor is below the dew point (301.2 K) and water phase humidity covers the whole surface. The hot air close to the floor that comes into contact with the cold surface, reduces its temperature and flows down due to the gravity. The remained hot air flows up to the roof.

In **Figure 10** the vertical temperature distribution in the region close to the floor calculated by three different numerical schemes (HYBRID, van LEER, SUPER) [6, 8, 12] is presented.

Temperature profile in the region of major gradient near the floor surface is predicted more abrupt by the SUPER scheme.

In **Figure 11** the vertical absolute humidity ratio (kg H₂O/kg of dry air) predicted by the three different numerical schemes (HYBRID, van LEER, SUPER) is presented at time 360 s.

Heat convection is accompanied by mass transfer and phase change of humid air. The larger gradient of temperature profile predicted by the SUPER scheme [10] leads to the formation of larger amount of water phase. Comparing the performance of the discretization schemes a more accurate solution of the condensation procedure is observed when applying the SUPER scheme.

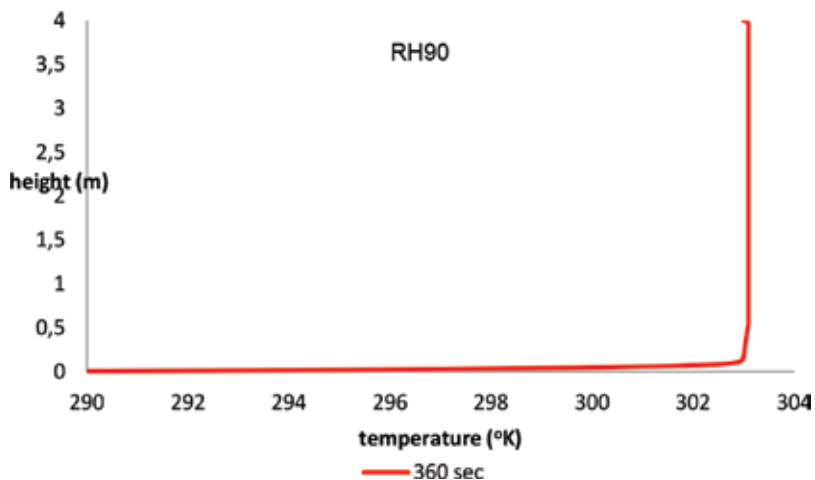


Figure 9. Vertical temperature distribution at time 360 s for initial humidity condition 90%.

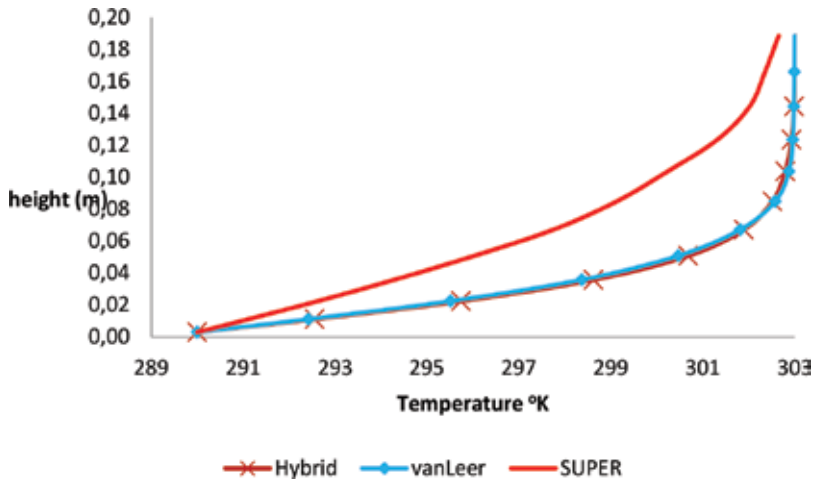


Figure 10. Vertical temperature distribution at time 360 s for initial humidity condition 90% in the middle of the domain at the height (0–0.2 m) calculated by three different numerical schemes.

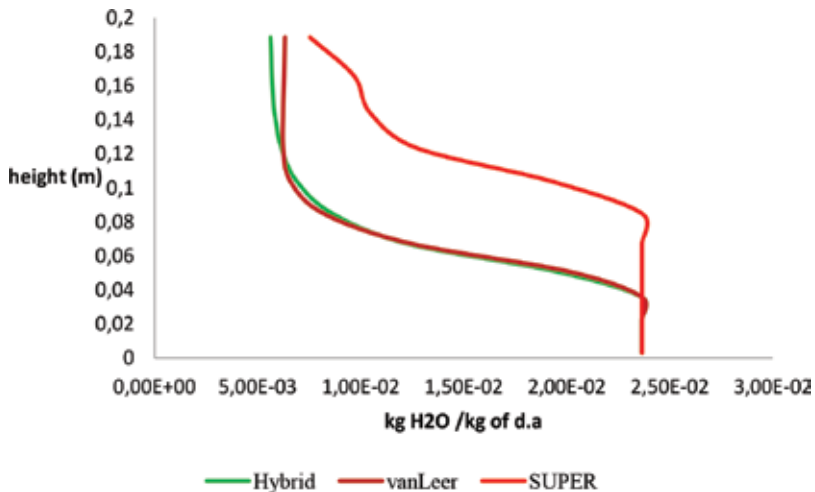


Figure 11. Vertical absolute humidity ratio (kg H₂O/kg of dry air) distribution at time 360 s for initial humidity condition 90% applying the van LEER, the HYBRID, the SUPER numerical schemes at height 0–0.2 m.

3. Conclusions

Numerical diffusion is the main source of errors in computational calculations. In this study the performance of various numerical schemes is evaluated in order to find the limitations of the convection formulations. The first-order UPWIND scheme presents more stability and overcomes the streamwise diffusion but such the higher-order schemes do not deal at all with the cross-stream diffusion. The second-order hybrid numerical schemes take into account the influence of more neighboring points but appears unstable at high Peclet numbers. The non-linear van LEER numerical scheme is highly stable and reduces satisfactory the false diffusion

errors. The flow oriented SUPER scheme overcomes the phenomenon of the numerical diffusion in most of the cases investigated without increasing the computational cost.

Acknowledgements

The first author (D. P. Karadimou) gratefully acknowledges the financial support from the State Scholarships Foundation of Greece through the "IKY Fellowships of Excellence for Postgraduate studies in Greece - SIEMENS" Program.

Author details

Despoina P. Karadimou^{1*} and Nikos-Christos Markatos^{1,2}

*Address all correspondence to: dkaradimou@gmail.com and n.markatos@ntua.gr

1 National Technical University of Athens, School of Chemical Engineering, Athens, Greece

2 Texas A & M University at Qatar, Education City, Doha, Qatar

References

- [1] Patel MK, Markatos NC, Cross M. Technical note-method of reducing false diffusion errors in convection diffusion problems. *Applied Mathematical Modelling*. 1985;**9**: 302-306
- [2] Patel MK, Markatos NC. An evaluation of eight discretization schemes for two-dimensional convection-diffusion equations. *International Journal for Numerical Methods in Fluids*. 1986;**6**:129-154
- [3] Darwish M, Moukalled F. A new approach for building bounded skew-upwind schemes. *Computer Methods in Applied Mechanics and Engineering*. 1996;**129**:221-233
- [4] De Vahl Davis G, Mallinson GD. An evaluation of upwind and central difference approximations by a study of recirculating flows. *Computers & Fluids*. 1976;**4**(1):29-43
- [5] Fromm JE. A method for reducing dispersion in convective difference schemes. *Journal of Computational Physics*. 1968;**3**:176-189
- [6] Spalding DB. A novel finite-difference formulation for different expressions involving both first and second derivatives. *International Journal for Numerical Methods in Engineering*. 1972;**4**:551-559
- [7] Leonard BP. A stable and accurate convective modelling procedure based on quadratic upstream interpolating. *Computational Mechanics and Applied Mechanical Engineering*. 1979;**4**:557-559

- [8] Van Leer B. Upwind-difference methods for aerodynamics problems governed by the Euler equations. Lectures in. Applications of Mathematics. 1985;22:327-336
- [9] Raithby GD. Skew upstream differencing schemes for problems involving fluid flow. Computer Methods in Applied Mechanics and Engineering. 1976;9:151-156
- [10] Patel MK, Markatos NC, Cross M. An assessment of flow oriented schemes for reducing 'false diffusion'. International Journal for Numerical Methods in Engineering. 1988;26:2279-2304
- [11] Carey C, Scanlon TJ, Fraser SM. SUCCA-an alternative scheme to reduce the effects of multidimensional false diffusion. Applied Mathematical Modelling. 1993;17(5):263-270
- [12] Karadimou DP, Markatos NC. A novel flow oriented discretization scheme for reducing false diffusion in three dimensional (3D) flows: An application in the indoor environment. Atmospheric Environment. 2012;61:327-339
- [13] Versteeg HK, Malalasekera W. An Introduction to Computational Fluid Dynamics-The Finite Volume Method. Essex, England: Longman Group Ltd; 1995
- [14] Spalding DB. Numerical computation of multiphase flow and heat-transfer. In: Taylor C, Morgan K, Michigan, United States, editors. Contribution to Recent Advances in Numerical Methods in Fluids. Pineridge Press; 1978. pp. 139-167
- [15] Markatos NC. Modelling of two-phase transient flow and combustion of granular propellants. International Journal of Multiphase Flow. 1983;12(6):913-933
- [16] Spalding DB. A general purpose computer program for multi-dimensional one or two-phase flow. Mathematics and Computers in Simulation XII. 1981:267-276
- [17] Gomma A, Halim MA, Elsaid AM. Experimental and numerical investigations of a triple-concentric tube heat exchanger. Applied Thermal Engineering. 2016;99:1303-1315
- [18] Benavides A, Van Wachem B. Eulerian-Eulerian prediction of dilute turbulent gas-particle flow in a backward facing step. International Journal of Heat and Fluid Flow. 2009;30:452-461
- [19] Chen F, Yu S, Lai A. Modeling particle distribution and deposition in indoor environments with a new drift-flux model. Atmospheric Environment. 2009;40:357-367
- [20] Yakhot V, Orszag SA, Thangam S, Gatski TB, Speziale CG. Development of turbulence models for shear flows by a double expansion technique. Physics of Fluids A. 1992;4(7):1510-1520
- [21] Markatos NC. The mathematical modelling of turbulent flows. Applied Mathematical Modelling. 1986;10:190-220
- [22] Padfield T. Conservation Physics. An Online Textbook in Serial Form. Available from: <http://www.conservaionphysics.org/atmcalc/atmoclc2.pdf>
- [23] Arnold Wexler, ASHRAE Handbook: Thermodynamic Properties of Dry Air, Moist Air and Water and S.I Psychrometric Charts, New York; 1983

Simulation of Natural Convection in Porous Media by Boundary Element Method

Janja Kramer Stajniko, Renata Jecl and Jure Ravnik

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.71230>

Abstract

In this chapter, the boundary element method (BEM) is introduced for solving problems of transport phenomena in porous media domains, which is an important topic in many engineering and scientific branches as well as in fields of practical interest. The main objective of the present work is to find a numerical solution of the governing set of equations written for fluid flow in porous media domains, representing conservation of mass, momentum, and energy. The momentum equation is based on the macroscopic Navier-Stokes equations and is coupled with the energy equation. In order to use BEM for the solution of the obtained set, the governing equations are transformed by the velocity-vorticity formulation, which separates the computational scheme into kinematic and kinetic computational parts. A combination of single- and sub-domain BEM is used to solve the obtained set of partial differential equations. Solution to a problem of natural convection in porous media saturated with pure fluid and nanofluid, respectively, for examples of 2D and 3D geometries, is shown. Results are compared to published work in order to estimate the accuracy of developed numerical algorithm. Based on the results, the applicability of the BEM for solving wide range of various problems is stated.

Keywords: boundary element method, porous media, velocity-vorticity formulation, natural convection, nanofluids

1. Introduction

Problems of transport phenomena in porous media have been widely investigated over the last few decades, mainly because of several important applications, which could be found in industry and environment, e.g., building insulation systems, dispersion of contaminants through water saturated soil, protection of groundwater resources, combustion technology.

Buoyancy driven flows in porous enclosures have been simulated using different mathematical models and numerical techniques. Most commonly used mathematical model of governing momentum equation is the Darcy's law, which is valid for the laminar flow regime ($Re < 10$), where the velocities are low and the viscous forces are predominant over inertia forces [1]. Extensions of the governing momentum equations have been made by analogy with the Navier-Stokes equations with addition of Brinkman term in order to consider the viscous diffusion and Forchheimer term to study the inertia effects on the free convection [2].

Problems of natural convection in porous media were studied intensively in last few decades, mainly for the cases of two-dimensional geometries. Two types of geometries are commonly investigated: porous enclosures where temperature gradient is imposed horizontally [3–7] or vertically [8–11]. Studies considering three-dimensional geometries are rare and are usually confined on using a simplified mathematical model, e.g., Darcy model or to conditions of heating from below [12–18]. Researches considering three-dimensional cavities with the condition of heating from the side were published in [19–21].

Recently, several researchers have been investigated buoyant flow in porous media domains saturated with nanofluids [22–24]. Nanoscale particles are often added to working fluids in order to enhance heat transfer or cooling processes. A comprehensive review of the studies considering convection heat transfer in porous media saturated with nanofluid was published in [25].

The solutions of the problems of transport phenomena in porous media have been obtained using different numerical methods, where the most commonly used methods are the finite element method (FEM), the finite difference method (FDM), and the finite volume method (FVM). As an alternative to others, in engineering practice widely used methods, the BEM was developed mainly because it was very efficient for solving potential problems of fluid mechanics (inviscid fluid flow, heat conduction, etc.), where the mathematical transformation of the governing set of partial differential equation results in boundary integral equations only. To rewrite the partial differential equation into an equivalent integral representation, the known fundamental solutions of the differential operator [26] and the Green's theorem are used. The discretized system contains only a fully populated system of integrals over boundary elements, which represent the main advantage over the volume-based methods.

When dealing with nonhomogenous and nonlinear problems, e.g., diffusion-convection problems, the domain integrals occur in the integral representation as well, which demands the extension of the classical BEM in order to additionally deal the problem within the domain. The main issue in this case is the evaluation of the domain matrices, which are full and unsymmetrical and require a lot of storage space. Several techniques have been developed in order to eliminate the domain integrals or transform them into the boundary integrals. One of the possibilities is the dual reciprocity boundary element method (DRBEM), which transforms domain integrals into a finite series of boundary integrals [27–29]. The nonhomogenous term is expanded in terms of radial basis functions. Since the discretization of the domain is represented only by grid points and the discretization of the geometry and fields on the boundary is piecewise polygonal, the DRBEM is still more flexible and efficient against other numerical methods, e.g., FDM.

Another possible extension of BEM is the boundary domain integral method (BDIM), which enables solving of strong nonlinear problems, where the domain-based effects are dominant, typical for the examples of the diffusion-convection problems [30–32]. The numerical algorithm solves the velocity-vorticity formulation of the Navier-Stokes equations, which separate the numerical scheme into the kinematic and kinetic computational parts. Consequently, the pressure is removed from the field functions conservation equations, and the calculation of the boundary pressure values is eliminated. Further advantage of BDIM is efficient dealing with boundary conditions on the solid boundaries in case of solving the vorticity equation. The vorticity is calculated explicitly from the kinematic computational part without using any approximate formulae. Using the subdomain technique [33], the problem of fully populated system matrices and corresponding memory requirements can be importantly reduced. A very stable and accurate numerical description of coupled diffusion-convection problems follows the use of Green's functions of the appropriate linear differential operators instead of upwinding schemes of different orders, as this is the case in other domain-type numerical techniques, which also eliminate the oscillations in the numerical solutions. A wavelet compression method for a single-domain BEM in 2D was introduced in [34].

In this chapter presented numerical algorithm is based on the combination of single- and subdomain BEM. The main advantage of the used method is that it enables an accurate prediction of the vorticity fields, which are in general defined as a curl of the velocity field. The vorticity is generated on the walls of the domain and influences the development of the flow field and furthermore the heat transfer. The single-domain BEM is used to solve the kinematics equation. The method is based on the fast multipole algorithm (FMM), which was introduced by Greengard and Rokhlin [35] for particle simulations and was later used for a wide variety of problems, e.g., for acceleration of the boundary integral Laplace equation by [36] and for coupling with BEM for the boundary matrices by [37]. The sub-domain BEM is used to solve the equations of the diffusion-advection type. A mesh of entire domain is made, where the integral equations are written for each of the sub-domains separately ([38–41]). Functions are discretized using the continuous quadratic boundary elements, whereas flux is discretized using discontinuous linear boundary elements, which enable to avoid flux definition problems in corners and edges. An over-determined system of linear equations is obtained, which is solved by a least squares manner.

A numerical approach based on the BEM has been used to solve a problem of buoyancy driven flows in porous media domain, saturated with pure fluid or nanofluid. The mathematical model is based on the Navier-Stokes equations, which are averaged over the representative elementary volume and rewritten into the velocity-vorticity formulation. The influence of several governing parameters, e.g., Rayleigh number, Darcy number, and volume fraction of nanoparticles, on the heat transfer and fluid flow characteristics is analyzed.

2. Mathematical model

2.1. Governing equations

The most general mathematical model for the transport phenomena in porous media is based on the volume-averaged Navier-Stokes equations, which are primarily written on the

microscopic level for the problems of pure fluid flow. Defining sufficiently large representative elementary volume (REV), with the restriction that only one part of the volume is available for the fluid flow, one can write the macroscopic set of governing equations. The REV is sufficiently large in case when it contains both, solid and fluid phases, irrespective of its position in porous media.

The presented governing equations are written for the case, when porous media are saturated with the nanofluid. The formulation enables considering both types of fluid flow by choosing correct parameter values. The properties describing nanofluids are density ρ_{nf} , dynamic viscosity μ_{nf} , heat capacitance $(c_p)_{nf}$, thermal expansion factor β_{nf} , and thermal conduction k_{nf} , where index nf stands for the nanofluid. Nanofluid properties are given in relation to pure fluid and pure solid properties linked with solid volume fraction of nanoparticles φ , which is given as:

$$\varphi = \frac{V_s}{V_s + V_f} \quad (1)$$

where V_s is volume of solid particles and V_f the volume of fluid. Relationships between nanofluid and pure fluid properties are described with models. A comprehensive review of different models can be found in [42]. In this chapter, macroscopic modeling of nanofluids is restricted to spherical nanoparticles, and it is suitable for small temperature gradients. Density of the nanofluid ρ_{nf} is given as:

$$\rho_{nf} = (1 - \varphi) \rho_f + \varphi \rho_s \quad (2)$$

where index f stands for the fluid phase and index s for the solid phase. The effective dynamic viscosity μ_{nf} can be given with according to [43] as:

$$\mu_{nf} = \frac{\mu_f}{(1 - \varphi)^{2.5\sigma}} \quad (3)$$

where the effective viscosity does not depend on the nanoparticle type. The heat capacitance of nanofluid can be given as:

$$(\rho c_p)_{nf} = (1 - \varphi) (\rho c_p)_f + \varphi (\rho c_p)_s \quad (4)$$

The nanofluid thermal expansion coefficient can be written in a similar way:

$$(\rho \beta)_{nf} = (1 - \varphi) (\rho \beta)_f + \varphi (\rho \beta)_s \quad (5)$$

taking into account the definition of ρ_{nf} it follows:

$$\beta_{nf} = \beta_f \left[\frac{1}{1 + \frac{\varphi \rho_s}{(1 - \varphi) \rho_f}} \frac{\beta_s}{\beta_f} + \frac{1}{1 + \frac{\varphi \rho_s}{(1 - \varphi) \rho_f}} \right] \quad (6)$$

The effective thermal conductivity k_{nf} is given with the Wasp model [44], which is valid only for the spherical particles, since it does not take into account the shape of the particles:

$$k_{nf} = k_f \frac{k_s + 2k_f - 2\phi(k_f - k_s)}{k_s + 2k_f + \phi(k_f - k_s)} \quad (7)$$

Further assumptions for the used model are: the nanoparticles are in thermal equilibrium with the base fluid and the nonslip boundary condition is considered. The fluid flow is assumed to be laminar, steady, Newtonian, and incompressible. The density depends only on the temperature variations and can be given with the Boussinesq approximation as:

$$\rho_{nf} = \rho_0(1 - \beta_{nf}(T - T_0)), \quad (8)$$

where T is temperature and index 0 refers to a reference state.

The problem of natural convection in saturated porous media can be described with the conservation equations for mass, momentum, and energy, written on the macroscopic level. The conservation of mass is given with the continuity equation:

$$\bar{\nabla} \cdot \bar{v} = 0 \quad (9)$$

where \bar{v} is volume averaged velocity vector.

The momentum equation is also known as Brinkman-Forchheimer equation and reads as:

$$\frac{1}{\phi} \frac{\partial \bar{v}}{\partial t} + \frac{1}{\phi^2} (\bar{v} \cdot \bar{\nabla}) \bar{v} = \frac{1}{\rho_{nf}} \bar{\nabla} p - \beta_{nf} (T - T_0) \bar{g} + \frac{1}{\phi} \frac{\mu_{nf}}{\rho_{nf}} \nabla^2 \bar{v} - \frac{1}{K} \frac{\mu_{nf}}{\rho_{nf}} \bar{v} - \frac{F \bar{v} |\bar{v}|}{K^{1/2}}, \quad (10)$$

where ϕ is porosity, t time, p pressure, T temperature, \bar{g} gravitational acceleration, K permeability, and F Forchheimer coefficient. There are two viscous and two inertial terms in the momentum equation. The third term on the r. h. s. of the equation is the Brinkman viscous term, which is analogous to the Laplacian term in the classical Navier-Stokes equations for pure fluid flow and expresses the viscous resistance or viscous drag force exerted by the solid phase on the flowing fluid at their contact surfaces. With the Brinkman term, the non-slip boundary condition on a surface which bounds porous media is satisfied [1]. The fourth term on the r. h. s. of the momentum equation is the Darcy term, where K is the permeability, which in general depends on the geometry of the porous medium and is a second-order tensor. When assuming an isotropic porous media, the permeability is a scalar. The last term in the momentum equation is the Forchheimer inertia term, which describes the nonlinear influences at higher velocities. The Forchheimer coefficient F is a dimensionless form-drag constant and is varying with the nature of porous medium. It can be written according to Ergun model as ([1]):

$$K = \frac{\phi^3 d_p^2}{a(1-\phi)^{2\prime}} F = \frac{b}{\sqrt{a} \phi^{3\prime}} \quad (11)$$

where a and b are Ergun's constants with values $a=150$ and $b=1.75$, while d_p is the average particle size of the bed.

Finally, the energy equation can be written as:

$$\sigma \frac{\partial T}{\partial t} + (\vec{v} \cdot \vec{\nabla})T = \frac{k_e}{(\rho c_p)_{nf}} \nabla^2 T, \quad (12)$$

where σ is the specific heat ratio, $\sigma = \phi + (1 - \phi)(\rho c_p)_p / (\rho c_p)_{nf}$, $(\rho c_p)_p$ and $(\rho c_p)_{nf}$ are heat capacities of solid and fluid phase, respectively. Furthermore, k_e is the effective conductivity of porous medium. It is assumed that the thermal properties of solid matrix and the nanofluid are identical [24, 45], resulting in $\sigma = 1$ and $k_e = k_{nf}$.

Governing equations (9), (10), and (12) can be converted into a nondimensional form by introduction of the following dimensionless variables:

$$\vec{v} \rightarrow \frac{\vec{v}}{v_0}, \quad \vec{r} \rightarrow \frac{\vec{r}}{L}, \quad t \rightarrow \frac{v_0 t}{L}, \quad \vec{g} \rightarrow \frac{\vec{g}}{g_0}, \quad p \rightarrow \frac{p}{p_0}, \quad T \rightarrow \frac{(T - T_0)}{\Delta T} \quad (13)$$

The parameters in the above expressions are v_0 characteristic velocity given as $v_0 = k_f / (\rho c_p)_f L$, which is common choice for buoyant flow simulations, k_f is the fluid thermal conductivity, $(\rho c_p)_f$ is the heat capacity for the fluid phase, and L is the characteristic length. Moreover, T_0 is characteristic temperature $T_0 = (T_2 - T_1)/2$ and ΔT is characteristic temperature difference $\Delta T = T_2 - T_1$, p_0 is the characteristic pressure $p_0 = 1\text{bar}$, while gravitational acceleration is $g_0 = 9.81 \text{ m/s}^2$.

In addition, the velocity-vorticity formulation of the governing equations is proposed by introduction of the vorticity vector, which is by the definition a curl of the velocity field $\vec{\omega} = \nabla \times \vec{v}$. The governing set of equations in nondimensional velocity-vorticity formulation can now be written in terms of kinematics equation, the vorticity transport equation, and the energy equation as:

$$\nabla^2 \vec{v} + \vec{\nabla} \times \vec{\omega} = 0 \quad (14)$$

$$(\vec{v} \cdot \vec{\nabla}) \vec{\omega} = (\vec{\omega} \cdot \vec{\nabla}) \vec{v} - C_A Pr Ra_T \phi^2 \vec{\nabla} \times T \vec{g} + C_B Pr \phi \nabla^2 \vec{\omega} - C_B \frac{Pr}{Da} \phi^2 \vec{\omega} - \frac{F}{Da} \phi^2 |\vec{v}| \vec{\omega}, \quad (15)$$

$$(\vec{v} \cdot \vec{\nabla})T = C_C \nabla^2 T \quad (16)$$

In the above equations, parameters C_A , C_B and C_C are presenting the nanofluid properties and are given with expressions:

$$C_A = \frac{\mu_{nf} \rho_f}{\mu_f \rho_{nf}}, \quad C_B = \frac{\beta_{nf}}{\beta_f}, \quad C_C = \frac{\alpha_{nf}}{\alpha_f} \quad (17)$$

where α_{nf} is thermal diffusivity of nanofluid, $\alpha_{nf} = k_{nf} / (\rho c_p)_{nf}$ and α_f thermal diffusivity of pure fluid $\alpha_f = k_f / (\rho c_p)_f$. The nanofluid properties are obtained using the expressions (2)–(8). For the

simulation of the pure fluid flow, the parameters are $C_A = C_B = C_C = 1$. The time derivatives in the vorticity and energy equations $\partial \omega / \partial t, \partial T / \partial t$ are omitted, since only steady flow simulations are shown in the present chapter.

The nondimensional parameters appearing in the momentum equation are:

- Fluid Rayleigh number $Ra_T = g \beta_T \Delta T L^3 \rho_f (\rho c_p)_f / \mu_f k_f$;
- Prandtl number $Pr = \mu_f c_p / k_f$;
- Darcy number $Da = K/L^2$;

The results are presented in terms of the porous Rayleigh number Ra_p , which links the thermal Rayleigh number and Darcy Number:

- $Ra_p = Ra_T \cdot Da$.

2.2. Boundary element method

The governing set of equations (14) and (15) in (16) is solved using an algorithm based on the combination of single-domain and sub-domain BEM, primarily developed for pure fluid flow simulations [39, 40] and later adopted for nanofluids [46] and porous media flow simulations [21]. The algorithm solves the velocity-vorticity formulation of Navier-Stokes equations. The sub-domain BEM solves the vorticity and energy transport equations. It is based on the domain decomposition, which results in sparse matrices and improves the efficiency of the solution to become comparable to FVM or FEM [47]. The kinematics equation for the calculation of the boundary vorticities is solved by a single-domain BEM. This results in full system of equations and limits the maximum grid size due to memory constraints. This drawback can be mitigated using the fast BEM, where sparse approximation of full matrices is used [37]. The main advantage of using the single-domain BEM for the boundary vorticity values is that the algorithm conserves mass in complex geometries, which is not the case when using velocity derivatives to calculate boundary vorticity values.

The numerical algorithm is devised as follows. At the beginning, the boundary conditions for the velocity and temperature are required and have to be given in terms of Dirichlet and Neumann type. In addition, the temperature and temperature flux on the solid walls and the no-slip boundary conditions are prescribed. The boundary conditions for the vorticity are unknown at the beginning and are calculated later as a part of numerical algorithm. The known boundary conditions are used to solve the kinematics equation (14) for the domain velocity values and energy equation (16) for the domain temperature values. The boundary vorticity values are first obtained using the single-domain BEM on the kinematics equation; moreover, the domain vorticity values are obtained out of vorticity transport equation (15) using a sub-domain BEM.

The outline of the numerical algorithm is as follows:

- Step 1.** Fluid/nanofluid and porous media properties are determined.
- Step 2.** Vorticity values on the boundary are calculated by a fast single-domain BEM from the kinematics equation (14).

Step 3. Velocity values within the domain are calculated by a sub-domain BEM from the kinematics equation (14).

Step 4. Temperature values are calculated by a subdomain BEM from the energy equation (16).

Step 5. Vorticity values within the domain are calculated by a subdomain BEM from the vorticity equation (17).

Step 6. Convergence check; all steps from 2 until 5 are repeated until all flow fields achieve the required accuracy.

In order to apply the proposed algorithm, governing equations have to be written in integral form. The integral representation is obtained using the Green's second identity for the unknown field function and for the fundamental solution of the Laplace equation as proposed in [48].

2.2.1. Integral representation of the kinematics equation

For the unknown boundary vorticity values, the single-domain BEM is used on the kinematics equation (14). The integral representation of the kinematics equation in its tangential form is:

$$c(\bar{\xi}) \bar{n}(\bar{\xi}) \times \bar{v}(\bar{\xi}) + \bar{n}(\bar{\xi}) \times \int_{\Gamma} \bar{v} \bar{\nabla} u^* \cdot \bar{n} d\Gamma = \bar{n}(\bar{\xi}) \times \int_{\Gamma} \bar{v} \times (\bar{n} \times \bar{\nabla}) u^* d\Gamma + \bar{n}(\bar{\xi}) \times \int_{\Omega} (\bar{\omega} \times \bar{\nabla} u^*) d\Omega \quad (18)$$

where Ω is the computational domain and $\Gamma = \partial\Omega$ is the boundary of the domain, $c(\bar{\xi})$ is geometric factor defined as $c(\bar{\xi}) = \theta/4\pi$, θ is the inner angle with origin in $\bar{\xi}$. If $\bar{\xi}$ lies inside the domain, then $c(\bar{\xi})=1$, and if $\bar{\xi}$ lies on a smooth boundary, then $c(\bar{\xi})=1/2$. Furthermore, \bar{n} is a vector normal to the boundary, and u^* is the fundamental solution of the Laplace equation given as:

$$u^* = \frac{1}{4\pi |\bar{\xi} - \bar{r}|} \quad (19)$$

The discretized system of equations is written for unknown boundary vorticities, while domain vorticity and velocity values are taken from the previous nonlinear iteration. The source point is set into every boundary node of the whole computational domain, which follows in full system matrix, where number of rows and columns are equal to number of boundary nodes. The system is solved using a LU decomposition method. The storage requirements are reduced with a kernel expansion approximation technique [49].

In addition, the kinematics equation (14) is used again in order to calculate domain velocity values with the sub-domain BEM. The following form of the integral equation is used:

$$c(\bar{\xi}) \bar{v}(\bar{\xi}) + \int_{\Gamma} \bar{v} (\bar{n} \cdot \bar{\nabla}) u^* d\Gamma = \int_{\Gamma} \bar{v} \times (\bar{n} \times \bar{\nabla}) u^* d\Gamma + \int_{\Omega} (\bar{\omega} \times \bar{\nabla} u^*) d\Omega. \quad (20)$$

The obtained integral kinematics equation is without the derivatives of the velocity or vorticity fields, which enables that the source point is set to function the nodes only. The domain

velocity values are calculated based on the known boundary values of the velocity from the initial boundary conditions, while the domain and boundary values of the vorticity are known from the previous iteration.

2.2.2. Integral representation of the vorticity and energy equations

In order to derive the integral form of the vorticity and energy equations, the same fundamental solution of the Laplace equation is used [26]. The final integral form of the vorticity transport equation is:

$$\begin{aligned}
 c(\xi)\omega_j(\xi) + \int_{\Gamma} \omega_j \bar{\nabla} u^* \cdot \bar{n} \, d\Gamma &= \int_{\Gamma} u^* q_j \, d\Gamma \\
 &+ \frac{1}{Pr} \frac{1}{C_B} \frac{1}{\phi} \int_{\Gamma} \bar{n} \cdot \left\{ u^* (\bar{\nabla} \omega_j - \bar{\omega} v_j) \right\} \, d\Gamma - \frac{1}{Pr} \frac{1}{C_B} \frac{1}{\phi} \int_{\Omega} (\bar{\nabla} \omega_j - \bar{\omega} v_j) \cdot \bar{\nabla} u^* \, d\Omega \\
 &- R a_T \frac{C_A}{C_B} \phi \int_{\Gamma} \left(u^* T \bar{g} \times \bar{n} \right)_j \, d\Gamma - R a_T \frac{C_A}{C_B} \phi \int_{\Gamma} \left(T \bar{\nabla} \times u^* \bar{g} \right)_j \, d\Omega \\
 &+ \frac{1}{Da} \phi \int_{\Omega} \omega_j u^* \, d\Omega + \frac{F}{Pr \sqrt{Da}} \frac{1}{C_B} \phi \left| \bar{\nabla} \right| \int_{\Omega} \omega_j u^* \, d\Omega,
 \end{aligned} \tag{21}$$

and finally the integral form of the energy transport equation reads as:

$$c(\xi)T(\xi) + \int_{\Gamma} T \bar{\nabla} u^* \cdot \bar{n} \, d\Gamma = \int_{\Gamma} u^* q_T \, d\Gamma + \frac{1}{C_C} \left[\int_{\Gamma} \bar{n} \cdot \{ u^* (\bar{\nabla} T) \} \, d\Gamma - \int_{\Omega} (\bar{\nabla} T) \cdot \bar{\nabla} u^* \, d\Omega \right] \tag{22}$$

In the above equations, q_j is a component of vorticity flux, whereas q_T is a heat flux. In the subdomain BEM method, a mesh of the entire domain Ω is made, each mesh element is named a subdomain. All equations are written for each of the subdomains. The field functions and flux across the boundary and within the domain are interpolated using shape functions. The hexahedral subdomains with 27 nodes are used, enabling continuous quadratic interpolation of field functions. The field functions on each element are interpolated using continuous quadratic interpolation, while fluxes are interpolated using the discontinuous linear interpolation. With discontinuous interpolation, the definition problems in corners and edges are avoided.

3. Test cases

The physical models where the above developed numerical scheme was tested are a two-dimensional rectangular enclosure and a three-dimensional cubical enclosure filled with fully saturated porous medium. Porous medium is assumed to be isotropic, homogenous, and in thermal equilibrium with the fluid phase. The simulation of fluid flow and heat transfer through porous media domain of pure fluid and nanofluid, respectively, is presented. Two opposite vertical walls are subjected to a temperature differences, while the rest of the walls is adiabatic and impermeable. Geometry with corresponding boundary conditions is shown in **Figure 1**. The boundary conditions for the current problem are:

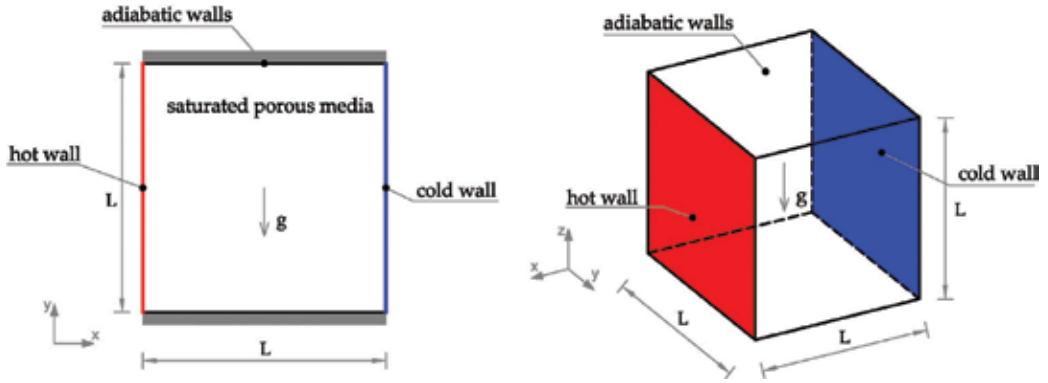


Figure 1. Two-dimensional and three-dimensional enclosures with corresponding boundary conditions.

$$\begin{aligned}
 v_x = v_y = 0, \omega = 0, T = T_H \quad \text{at } x = 0, \\
 v_x = v_y = 0, \omega = 0, T = T_C \quad \text{at } x = L, \\
 v_x = v_y = 0, \omega = 0, \frac{\partial T}{\partial y} = 0 \quad \text{at } y = 0 \text{ and } z = 0, \\
 v_x = v_y = 0, \omega = 0, \frac{\partial T}{\partial y} = 0 \quad \text{at } y = L \text{ and } z = L
 \end{aligned} \tag{23}$$

Due to applied temperature gradient, density differences are induced, which result in appearance of thermal buoyancy force producing a large vortex in the main part of the cavity.

The overall heat transfer through porous media is expected to depend on several fluid and porous media properties, such as porosity, permeability, thermal conductivity, heat capacitance, solid volume fraction of nanoparticles, and types of nanoparticles. In order to compare different conditions on the heat transfer characteristics, the wall heat flux is calculated, which is expressed in terms of the dimensionless Nusselt number as:

$$Nu = \frac{k_{nf}}{k_f} \int_{\Gamma} \bar{\nabla} T \cdot \bar{n} d\Gamma, \tag{24}$$

where Γ is the surface through which the heat flux is calculated and \bar{n} is the unit normal to this surface. The definition is valid for nanofluids as well as for pure fluids, since there the ratio of thermal conductivities is $k_{nf}/k_f = 1$.

In the present study, the Cu nanoparticles are added to the water as a base fluid. The thermo-physical properties of Cu nanoparticles and water are given in **Table 1** [50].

In order to obtain a grid-independent solution, at the beginning, a grid sensitivity analysis was performed. Two nonuniform meshes for 2D geometry and four nonuniform meshes for 3D geometry have been tested. The results are shown for the case, when porous media are saturated with pure fluid and parameters $Ra_p = 100$, $Pr = 0.71$, $\phi = 0.8$, $\sigma = 1$ and various values of Da . The results are presented in **Table 2**.

	c_p [J/kg K]	ρ [kg/m ³]	K [W/m K]	β [$\times 10^{-5}$ K ⁻¹]	α [$\times 10^{-7}$ m ² /s]
Water	4179	997.1	0.613	21	1.47
Cu	385	8933	400	1.67	1163

Table 1. Thermophysical properties of water and Cu nanoparticles.

Mesh	Number of nodes	$Ra_p = 100, Pr = 0.71, \phi = 0.8$					
		Da	10^{-1}	10^{-2}	10^{-3}	10^{-4}	10^{-5}
2D	20×20	1681	1.0639	1.6329	2.3697	2.8756	3.1656
	30×30	3721	1.0638	1.6331	2.3680	2.8537	3.1503
3D	$12 \times 12 \times 12$	15,625	1.0423	1.5428	2.3432	2.9784	3.3008
	$20 \times 8 \times 20$	28,577	1.0394	1.5329	2.3313	2.9575	3.2950
	$22 \times 10 \times 22$	42,525	1.0393	1.5327	2.3307	2.9552	3.2945
	$30 \times 10 \times 30$	78,141	1.0393	1.5325	2.3303	2.9541	3.2934

Table 2. Variations of Nusselt number with different grid sizes and various Darcy numbers.

When comparing 2D and 3D results, it can be observed that for the case of low values of Darcy number, 2D simulation underestimates the heat transfer up to 4.5%. Based on the presented results, the 20×20 mesh for 2D simulations and $20 \times 8 \times 20$ for 3D simulations were chosen.

The validation of numerical code has been primarily performed for the pure fluid saturating the porous media. The results for 2D geometry, compared to [24, 51], are presented in **Table 3**. Results for 3D geometry are compared to [20] and are presented in **Table 4**. Furthermore, in **Table 5**, the results for natural convection in 2D enclosure for a nanofluid saturated porous media are shown and compared to [24]. According to the comparable study, the Cu nanoparticles were added to the water as a base fluid.

From the comparison, it can be observed that the results agree well with the data from the published studies, which confirm accuracy of the obtained numerical algorithm.

The isotherms for Cu-water nanofluid under different values of porous Rayleigh number and Darcy number at porosity $\phi = 0.4$ and different values of solid volume fraction are shown in **Figure 2**. Solid lines correspond to $\varphi = 0.0$, dotted lines to $\varphi = 0.025$, and dashed lines to $\varphi = 0.05$. Heat transfer in porous medium is mostly affected by a Rayleigh and Darcy numbers. At $Ra_p = 10$, the heat transfer in horizontal direction is weak, and the main heat transfer mechanism in this case is conduction. Increase of Ra_p results in stronger convective motion, which is clearly evident from the temperature field; when $Ra_p = 1000$, thin boundary layers are created near the hot and cold walls, and the isotherms in the core region become

Da	Ra _p	[51]	[24]	Present
10 ⁻²	10	1.015	1.010	1.012
	100	1.530	1.533	1.503
	1000	3.555	3.602	3.499
10 ⁻⁴	10	1.071	1.065	1.070
	100	2.725	2.764	2.777
	1000	9.183	9.454	9.174
10 ⁻⁶	10	1.079	1.072	1.093
	100	2.997	2.980	3.241
	1000	11.790	11.924	12.895

Table 3. Validation of the numerical code by a comparison of average Nu for natural convection in porous media saturated with a pure fluid for Pr = 1.0, φ=0.6 and different Da and Ra_p, for 2D geometry.

Ra _p = 1000, Pr = 0.71, φ = 0.8					
Da	10 ⁻²	10 ⁻³	10 ⁻⁴	10 ⁻⁵	10 ⁻⁶
[20]	3.99	6.95	10.14	12.78	13.72
Present	3.770	6.922	10.558	13.242	14.568

Table 4. Nusselt number values for the 3D natural convection in a cubic enclosure filled with porous media saturated with pure fluid.

Da	Ra _p	[24]	Present	[24]	Present	[24]	Present
φ = 0.05							
		φ = 0.4		φ = 0.6		φ = 0.9	
10 ⁻²	1000	3.433	3.400	3.850	3.826	4.162	4.145
10 ⁻⁴	1000	9.117	9.132	9.590	9.743	9.901	10.154
10 ⁻⁶	1000	11.778	12.991	11.899	13.128	11.976	13.195
φ = 0.4							
		φ = 0.0		φ = 0.025		φ = 0.05	
10 ⁻²	10	1.007	1.008	1.081	1.083	1.160	1.162
10 ⁻²	1000	3.302	3.282	3.370	3.345	3.433	3.400
10 ⁻⁶	1000	11.867	13.238	11.847	13.131	11.778	12.991

Table 5. Nusselt number values for a natural convection in porous media saturated with nanofluid in 2D enclosure for various governing parameters (Pr = 6.2).

almost horizontal and parallel to adiabatic and impermeable walls. According to the temperature fields, decrease of *Da* enhances the heat transfer through cavity. The *Da* number influences the magnitude of the Darcy term in the vorticity equation (10). With increase of

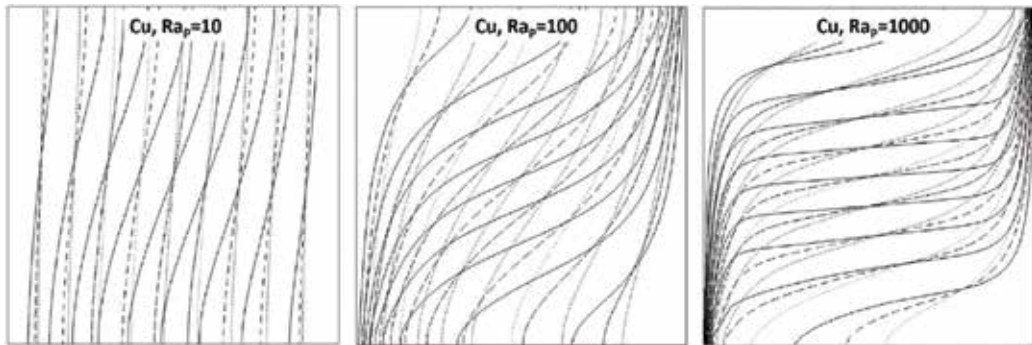


Figure 2. Temperature fields for different values of Ra_p for Cu nanofluid and different solid volume fractions for $Pr = 6.2$, $Da = 10^{-6}$; solid lines are for $\phi = 0.0$, dashed lines for $\phi = 0.025$ and dotted lines for $\phi = 0.05$.

Da , the flow regime is transited into the Darcy flow regime, which can be described by the Darcy’s law.

According to the results, the addition of nanoparticles into the water causes the attenuation of the convective motion. However, the overall heat transfer is enhanced with increase of solid volume fraction of nanoparticles in cases of conduction-dominated regimes (low values of Rayleigh numbers $Ra_p < 100$ and high values of Darcy numbers $Da > 10^{-4}$). On the other hand in convection dominated regimes ($Ra_p > 100$, $Da < 10^{-4}$), the addition of nanoparticles diminishes the convection, which results in lower values of Nusselt numbers. **Figure 3** shows the dependence of Nu on porous Rayleigh number and solid volume fraction of nanoparticles. It can be observed that for any values of Ra_p with increase of ϕ , the heat transfer increases.

When observing the flow structure in 3D enclosure, it is obvious that the flow field is not far from being 2D, which is a consequence of the applied temperature difference between the opposite walls, which causes large two-dimensional vortex in the y plane. In order to observe

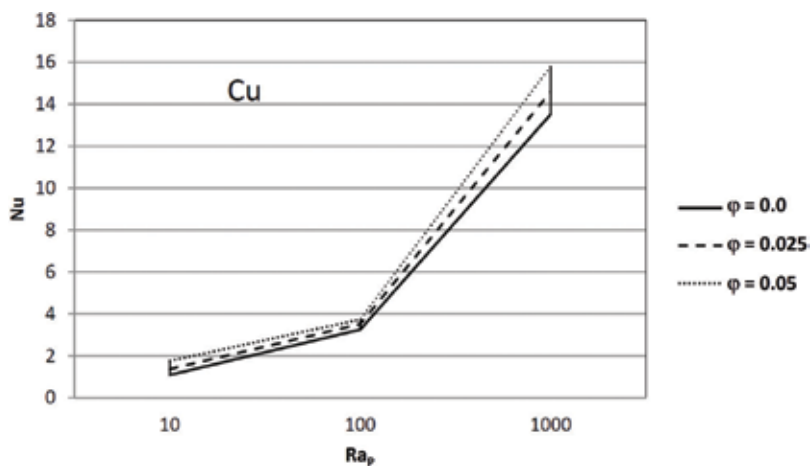


Figure 3. Nusselt number values depending on porous Rayleigh number for $Pr = 6.2$, $Da = 10^{-6}$ and different values of solid volume fraction of nanoparticles.

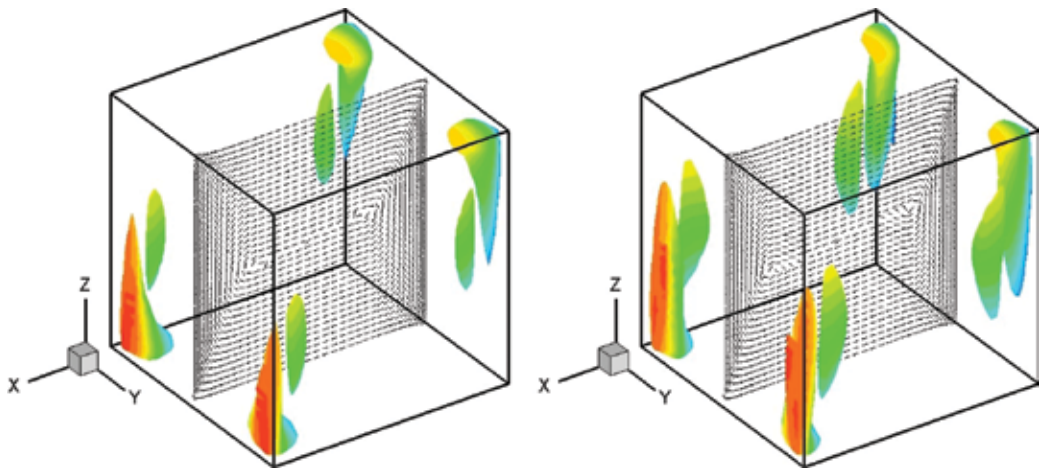


Figure 4. Iso-surfaces for $Ra_p=500$, $Da=10^{-3}$ and absolute value of velocity component $|v_y|=3$ (left) and $Ra_p=1000$, $Da=10^{-3}$, $|v_y|=7$ (right). Contours of temperature are displayed on the iso-surfaces ($-0.5 < T < 0.5$). In addition, the velocity vectors on the plane $y=0.5$ are displayed.

the 3D nature of the phenomena, the iso-surfaces of the absolute value of y velocity component are shown in **Figure 4**. The extent of movement perpendicular to the plane of the main vortex is small, but it becomes more apparent in case of higher values of Ra_p and lower values of Da .

4. Conclusion

The numerical method based on the BEM is presented for solving the coupled set of partial differential equations, which describe the fluid flow and heat transfer in porous medium domain. The mathematical model is based on the Navier-Stokes equations, which are averaged over the representative elementary volume. The proposed numerical algorithm solves the velocity-vorticity formulation of the governing equations. The numerical scheme is split into the single-domain BEM, which solves the kinematic equation for unknown boundary vorticity values and sub-domain BEM for domain velocity, vorticity, and temperature values.

The numerical algorithm is tested on an example of natural convection phenomena in porous media domain for a 2D as well as 3D geometry. Porous media are fully saturated with pure fluid or water-based nanofluid with addition of Cu nanoparticles. Obtained numerical results were validated with available benchmark solutions.

The natural convection phenomena strongly depend on the parameters, e.g., Rayleigh number and Darcy number. Addition of nanoparticles to a base fluid enhances the heat transfer through porous media, when conduction is the dominant heat transfer mechanism. On the other hand, in convection dominated regime, the addition of nanoparticles reduces the magnitude of convective motion.

The good agreement of the results with the published ones confirms the efficiency of the BEM-based methods as a powerful alternative to the existing numerical methods.

Author details

Janja Kramer Stajnik^{1*}, Renata Jecl¹ and Jure Ravnik²

Address all correspondence to: janja.kramer@um.si

1 Faculty of Civil Engineering, Transportation Engineering and Architecture, University of Maribor, Maribor, Slovenia

2 Faculty of Mechanical Engineering, University of Maribor, Maribor, Slovenia

References

- [1] Nield DA, Bejan A. Convection in Porous Media. 4th ed. United States of America: Springer; 2013
- [2] Lauriat G, Prasad V. Natural convection in a vertical porous cavity: A numerical study for Brinkman-extended Darcy formulation. *Journal of Heat Transfer*. 1987;**109**:688-696
- [3] Beckermann C, Viskanta R, Ramadhyani S. A numerical study of non-Darcian natural convection in a vertical enclosure filled with a porous medium. *Numerical Heat Transfer*. 1986;**10**:557-570
- [4] Lauriat G, Prasad V. Non-Darcian effects on natural convection in a vertical porous enclosure. *International Journal of Heat and Mass Transfer*. 1989;**32**:2135-2148
- [5] Jecl R, Škerget L, Petrešin E. Boundary domain integral method for transport phenomena in porous media. *International Journal for Numerical Methods in Fluids*. 2001;**35**:39-54
- [6] Baytas AC, Pop I. Free convection in a square porous cavity using a thermal nonequilibrium model. *International Journal of Thermal Sciences*. 2002;**41**:861-870
- [7] Basak T, Roy S, Paul T, Pop I. Natural convection in a square cavity filled with a porous medium: Effects of various thermal boundary conditions. *International Journal of Heat and Mass Transfer*. 2006;**49**:1430-1441
- [8] Prasad V, Kulacki FA. Natural convection in horizontal porous layers with localized heating from below. *Journal of Heat Transfer*. 1987;**109**:795-798
- [9] Kladias N, Prasad V. Natural convection in horizontal porous layers: Effects of Darcy and Prandtl numbers. *Journal of Heat Transfer*. 1989;**111**:926-935
- [10] Rosenberg ND, Spera FJ. Thermohaline convection in a porous medium heated from below. *International Journal of Heat and Mass Transfer*. 1992;**35**:1261-1273
- [11] Storesletten L. Natural convection in a horizontal porous layer with anisotropic thermal diffusivity. *Transport in Porous Media*. 1993;**12**:19-29
- [12] Beck JL. Convection in a box of porous material saturated with fluid. *Physics of Fluids*. 1972;**15**:1377-1383

- [13] Holst PH, Aziz K. Transient three-dimensional natural convection in confined porous media. *International Journal of Heat and Mass Transfer*. 1972;**15**:73-90
- [14] Zebib a, Kassoy DR. Onset of natural convection in a box of water-saturated porous media with large temperature variation. *The Physics of Fluids*. 1977;**20**:4-9
- [15] Zebib A, Kassoy DR. Three-dimensional natural convection motion in a confined porous medium. *Physics of Fluids*. 1978;**21**:1-3
- [16] Horne RN. Three-dimensional natural convection in a confined porous medium heated from below. *Journal of Fluid Dynamics*. 1979;**92**:25-38
- [17] Neto HL, Quaresma JNN. Natural convection in three-dimensional porous cavities: Integral transform method. *International Journal of Heat and Mass Transfer*. 2002; **45**:3013-3032
- [18] Sezai I. Flow patterns in a fluid-saturated porous cube heated from below. *Journal of Fluid Mechanics*. 2005;**523**:393-410
- [19] Dawood AS, Burns PJ. Steady three-dimensional convective heat transfer in a porous box via multigrid. *Numerical Heat Transfer, Part A*. 1992;**22**:167-198
- [20] Sharma RV, Sharma RP. Non-Darcy effects on three-dimensional natural convection in a porous box. *Annals of the Assembly for International Heat Transfer Conference*. 2006
- [21] Kramer J, Ravnik J, Jecl R, Škerget L. Simulation of 3D flow in porous media by boundary element method. *Engineering Analysis with Boundary Elements*. 2011;**35**:1256-1264
- [22] Shermet MA, Pop I. Conjugate natural convection in a square porous cavity filled by a nanofluid using Buongiorno's mathematical model. *International Journal of Heat and Mass Transfer*. 2014;**79**:137-145
- [23] Grosan T, Revnic C, Pop I, Ingham DB. Free convection heat transfer in a square cavity filled with a porous medium saturated by a nanofluid. *International Journal of Heat and Mass Transfer*. 2015;**87**:36-41
- [24] Nguyen MT, Aly AM, Lee SW. Natural convection in a non-Darcy porous cavity filled with Cu-water nanofluid using the characteristic-based split procedure in finite-element method. *Numerical Heat Transfer, Part A*. 2015;**67**:224-247
- [25] Mahdi RA, Mohammed HA, Munisamy KM, Saeid NH. Review of convection heat transfer and fluid flow in porous media with nanofluid. *Renewable and Sustainable Energy Reviews*. 2015;**41**:715-734
- [26] Wrobel LC. *The Boundary Element Method*. Chichester, England, New York: John Wiley & Sons Ltd; 2002
- [27] Partridge P, Brebbia C, Wrobel L. *The Dual Reciprocity Boundary Element Method*. Southampton: Computational Mechanics Publications; 1992
- [28] Perez-Gavilan JJ, Aliabadi MH. A Galerkin boundary element formulation with dual reciprocity for elastodynamics. *International Journal of Numerical Methods in Engineering*. 2000;**48**:1331-1344

- [29] Blobner J, Hriberšek M, Kuhn G. Dual reciprocity BEM-BDIM technique for conjugate heat transfer computations. *Computer Methods in Applied Mechanics and Engineering*. 2000;**190**:1105-1116
- [30] Škerget L, Hriberšek M, Kuhn G. Computational fluid dynamics by boundary domain integral method. *International Journal for Numerical Methods in Engineering*. 1999; **46**:1291-1311
- [31] Jecl R, Škerget L, Petrešin E. Boundary domain integral method for transport phenomena in porous media. *International Journal for Numerical Methods in Engineering*. 2001;**35**:39-54
- [32] Kramer J, Jecl R, Škerget L. Boundary domain integral method for study of double diffusive natural convection in porous media. *Engineering Analysis with Boundary Elements*. 2007;**31**:897-905
- [33] Hriberšek M, Škerget L. Iterative methods in solving Navier-stokes equations by the boundary element method. *International Journal for Numerical Methods in Engineering*. 1996;**39**:115-139
- [34] Ravnik J, Škerget L, Hriberšek M. The wavelet transform for BEM computational fluid dynamics. *Engineering Analysis with Boundary Elements*. 2004;**28**:1303-1314
- [35] Greengard L, Rokhlin V. A fast algorithm for particle simulations. *Journal of Computational Physics*. 1997;**135**:280-292
- [36] Popov V, Power H, Walker SP. Numerical comparison between two possible multipole alternatives for the BEM solution of 3D elasticity problems based upon Taylor series expansions. *Engineering Analysis with Boundary Elements*. 2003;**27**:521-531
- [37] Ravnik J, Škerget L, Žunič Z. Comparison between wavelet and fast multipole data sparse approximations for Poisson and kinematic boundary-domain integral equations. *Computer Methods in Applied Mechanics and Engineering*. 2009;**198**:1473-1485
- [38] Ramšak M, Škerget L. 3D multidomain BEM for solving the Laplace equation. *Engineering Analysis with Boundary Elements*. 2007;**31**:528-538
- [39] Ravnik J, Škerget L, Žunič Z. Velocity-vorticity formulation for 3D natural convection in an inclined enclosure by BEM. *International Journal of Heat and Mass Transfer*. 2008;**51**:4517-4527
- [40] Ravnik J, Žunič Z. Combined single domain and subdomain BEM for 3D laminar viscous flow. *Engineering Analysis with Boundary Elements*. 2009;**33**:420-424
- [41] Popov V, Power H, Škerget L, editors. *Domain Decomposition Techniques for Boundary Elements: Applications to Fluid Flow*. Southampton, Billerica: WIT Press; 2007
- [42] Haddad Z, Oztop HF, Abu-Nada E, Mataoui A. A review on natural convective heat transfer of nanofluids. *Renewable and Sustainable Energy Reviews*. 2012;**16**:5363-5378
- [43] Brinkman HC. The viscosity of concentrated suspensions and solutions. *The Journal of Chemical Physics*. 1952;**20**:571-581

- [44] Wasp FJ. Solid-liquid slurry pipeline transportation. Clausthal Trans Tech Publications; 1977
- [45] Bourantas GC, Skouras ED, Loukopoulos VC, Burganos VN. Heat transfer and natural convection of nanofluids in porous media. *European Journal of Mechanics B/Fluids*. 2014;**43**:45-56
- [46] Ravnik J, Hriberšek M. Analysis of three-dimensional natural convection of nanofluids by BEM. *Engineering Analysis with Boundary Elements*. 2010;**34**:1018-1030
- [47] Ravnik J, Hriberšek M. 2D velocity vorticity based LES for the solution of natural convection in a differentially heated enclosure by wavelet transform based BEM and FEM. *Engineering Analysis with Boundary Elements*. 2006;**30**:671-686
- [48] Škerget L, Žunič Z. Natural convection flows in complex cavities by BEM. *International Journal of Numerical Methods for Heat & Fluid Flow*. 2003;**13**:720-735
- [49] Ravnik J, Žunič Z. Fast single domain-subdomain BEM algorithm for 3D incompressible fluid flow and heat transfer. *International Journal for Numerical Methods in Engineering*. 2009;**77**:1627-1645
- [50] Oztop HF, Abu-Nada E. Natural convection of water based nanofluids in an inclined enclosure with heat source. *International Journal of Heat and Fluid Flow*. 2008;**29**:1326-1336
- [51] Nithiarasu P, Seetharamu KN, Sundararajan T. Natural convective heat transfer in a fluid saturated variable porosity medium. *International Journal of Heat and Mass Transfer*. 1997;**40**:3955-3967

Numerical Simulations of a High-Resolution RANS-FVDM Scheme for the Design of a Gas Turbine Centrifugal Compressor

Chellapilla V K N S N Moorthy and Vadapalli Srinivas

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.72098>

Abstract

The aero-thermodynamic design and performance of a compressor need to conquer many vital challenges like it is a gas-driven turbo-machinery component, involvement of extensive iterative process for the convergence of the design, enormous design complexity due to three-dimensional flow phenomena, and multiflow physics embedded within a dynamic state-of-the-art. In this chapter, a strong attempt is made to address the above-cited technical issues to achieve an optimized design and performance of a centrifugal compressor with backward swept blade profile producing total pressure ratio of 5.4 with an ingested mass flow rate of 5.73 kg/s. A mean-line design methodology was implemented to configure sizing of the compressor. An optimum grid size was well validated by carrying out computational analysis with three different mesh sizes within the same framework. Finally, a detailed three-dimensional numerical simulation was performed using Reynolds-averaged Navier-Stokes equations based on finite volume discretization method (RANS-FVDM) scheme. Consequently, the polytropic efficiency, total-to-total efficiency, stagnation pressure ratio at a fixed rotational speed, and the overall design and aero-thermodynamic performance of the centrifugal compressor are validated.

Keywords: numerical simulations, Reynolds-averaged Navier-Stokes equations, finite volume discretization method, compressor design

1. Introduction

Centrifugal compressors find usage over wide range of propulsion applications and are regarded as one among the key air-breathing propulsive engine components. The cognitive research and development of compressors is directed toward achieving a higher pressure ratio, higher

efficiency, and reduced structural weight of compressor and the engine as well. Various compressor stages achieve gradual increase in the stagnation-to-flow pressure contributed by flow diffusion. Energy is added in the rotor blade section, increasing the total pressure and absolute component of flow velocity. Stator blade row diffuses the flow, thus reducing absolute velocity component and elevating static pressure. Blade topology requires adaptation of a cautious design procedure to achieve the designated pressure rise while minimizing aero-thermodynamic losses in order to run and achieve design pressure ratios and design efficiencies.

In this chapter, a strong attempt is made to enumerate the detailed procedure of the centrifugal compressor in Section 2, concepts and basics of Numerical Schemes in Section 3. Further, a case study is discussed in detail in Section 4, the corresponding results and discussions are presented in Section 5. Final conclusions are summarized in Section 6.

2. Centrifugal compressor and the basic steps involved in its design

The design of the centrifugal compressor is very difficult due to the reasons like it is a gas-driven turbo-machinery component, involvement of extensive iterative process for the convergence of the design, enormous design complexity due to three-dimensional flow phenomena, and multiflow physics embedded within a dynamic state-of-the-art. Hence, in this article, a detailed procedure for the design is presented based on the design methodology as explained in the literature [1–3]. The basic steps involved in the design of the centrifugal compressor are discussed in detail in the following sections.

2.1. Impeller and velocity triangles

Impeller is the main design task during the phase of the compressor design. There are various strategic geometric/design features to be identified and discussed with respect to impeller, its inlet and outlet design. In principal, aerodynamic losses occurring in majority of turbo-machines arise primarily due to the boundary layer growth, its separation on blade profile, and passage surfaces termed as profile losses widely configured under primary losses.

The velocity triangles of the impeller of a centrifugal compressor are diagrammatically shown in **Figure 1**.

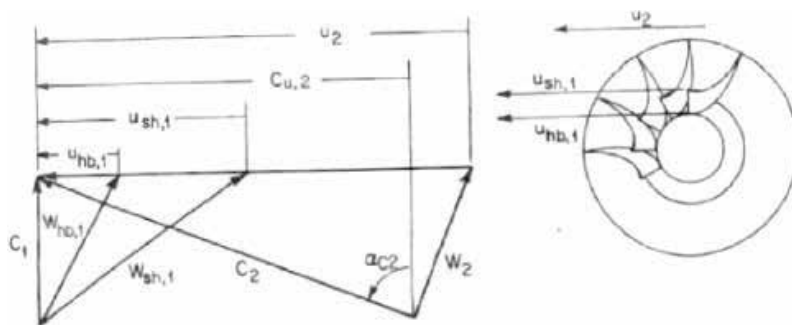


Figure 1. Velocity triangles of impeller for a centrifugal compressor.

The impeller design is the basic and major vital part of the design of the compressor. The following conditions for the design may be considered as limiting conditions

1. At the inlet of the impeller, very high relative velocity of the shroud ($W_{sh,1}$) should be avoided.
2. Separation of the flow in the passages of the impeller should be minimized so that the losses are minimized.
3. To meet the above limiting condition, the de Haller impeller ratio should be $\frac{W_2}{W_{sh,1}} > 0.75$.
4. The stability of the diffuser strongly depends on $\alpha_{c,2}$ and for both vaneless and vaned diffusers, $\alpha_{c,2} \leq 70^\circ$. Further, it should be taken from 60 to 65° for better design.
5. For the backswept impellers, the ratio of the relative velocities should be increased for betterment in the design.
6. Little bit of the increase in the stagnation enthalpy will be attributed to that of the velocity, because of which the polytropic efficiency of impellers with nil back sweep will be less than that of backswept impellers.
7. Work coefficient for radial impellers with zero back sweep is around 0.9 and 0.5 for backswept impellers.
8. To with stand the stress levels, the maximum rotational speed of the backswept impeller is maintained less than that of the radial impeller with same tip diameter.
9. The above two conditions lead to decrease in the maximum possible work per stage for a backswept impellers.
10. Hence, if the maximum work per stage is not required, blade swept angles up to 50° can be used.

2.2. Specific work

The specific work can be calculated using the following correlation, which is based on energy and Euler turbine equations

$$(u_2 C_{u,2} - u_1 C_{u,1}) = g_c C_p g T_{0,1} \left\{ \left(\frac{p_{0,2}}{p_{0,1}} \right) \left[\left(\frac{R}{C_p} \right)^{\frac{1}{\eta_{p,c}}} \right] - 1 \right\} \quad (1)$$

2.3. Slip factor

In general for the design of the compressor, the inlet swirl is considered as zero and hence, $C_{u,1} = 0$. Furthermore, the Wiesner's correlation for slip factor is given by

$$\sigma_w = slipfactor = \frac{C_{u,2,ac}}{C_{u,2,tl}} = 1 - \frac{\sqrt{\cos \beta_2}}{Z^{0.7}} \quad (2)$$

where number of blades is denoted by Z and β_2 represents the angle between the radial direction and tangent to the rotor blade at the periphery.

$\alpha_{c,2}$ should be less than 60° such that the downstream diffuser does not get prone to stall when C_{2c} and $C_{r,2}$ are reduced by keeping very large $\alpha_{c,2}$.

2.4. Number of blades

The number of blades and the exit blade angle are dependent on each other and can be calculated using loading coefficient. Loading coefficient ψ is the ratio between the outlet tangential flow velocity and blade speed, which is given by the correlation

$$\psi = \frac{C_{u,2,sc}}{u_2} = \left\{ \left[\frac{\tan\beta_2}{\tan\alpha_{c,2}} \right] + \frac{1}{\sigma_w} \right\}^{-1} \quad (3)$$

For the sake of easy reference, even number of blades are chosen so that half of the blades can be considered as splitter blades. In general for a better design, 20 blades and 45° blade angle at outlet are taken.

2.5. Blade peripheral speed

The peripheral speed of the impeller is calculated using the following correlations

$$-\Delta h_0 = \psi u^2 \quad (4)$$

where increase in the enthalpy is given by

$$\Delta h_0 = C_p(T_{0,2} - T_{0,1}) \quad (5)$$

The isentropic law can be applied, which is given by

$$\frac{T_{0,2}}{T_{0,1}} = \left(\frac{p_{0,2}}{p_{0,1}} \right)^{\left[\left(\frac{\gamma}{c_p} \right) \frac{1}{\eta_{p,c}} \right]} \quad (6)$$

The polytropic efficiency $\eta_{p,c}$ can be calculated by choosing isentropic efficiency $\eta_{s,c}$ as 0.83 and optimum specific speed N_s as 0.8 for the better design and using the following correlation

$$\eta_{s,c} = \frac{r^{R/c_p - 1}}{r^{(R/c_p)\eta_{p,c} - 1}} \quad (7)$$

Other blade parameters can be calculated using the regular trigonometric correlations from the velocity triangles of the blade like the following

$$C_{u,2,ac} = \psi u_2 \quad (8)$$

$$C_{r,2} = \frac{C_{u,2,ac}}{\tan \alpha_{c,2}} \quad (9)$$

$$C_2 = \frac{C_{u,2,ac}}{\sin \alpha_{c,2}} \quad (10)$$

2.6. Density at the inlet of the blade

The Mach number at the inlet can be calculated using the following correlation

$$\frac{C}{\sqrt{g_c R T_0}} = \sqrt{2 \frac{C_p}{R} \left[1 - \left[1 + \frac{M^2}{2 \left(\frac{C_p}{R} - 1 \right)} \right]^{-1} \right]} \quad (11)$$

Static density can be calculated from the following correlations

$$\frac{\rho_0}{\rho_{st}} = \left[1 + \frac{M^2}{2 \left(\frac{C_p}{R} - 1 \right)} \right]^{\left(\frac{C_p}{R} - 1 \right)} \quad (12)$$

2.7. Rotational speed

The rotational speed can be calculated using the following correlation

$$N = \frac{60 N_s (g_c \Delta h_0)^{3/4}}{2\pi \sqrt{\dot{V}}} \quad (13)$$

2.8. Blade axial width at outlet

The outlet blade axial width is calculated with the assumptions that the effect of the thickness and boundary layer of the blade are neglected. The following correlations are used for the calculations

$$\dot{m} = C_{r,2} \rho_{st,2} \pi d_2 b_2 \quad (14)$$

where, d_2 , the rotor diameter can be calculated using the correlation

$$u_2 = \frac{\pi d_2 N}{60} \quad (15)$$

2.9. Flow separation condition

Flow separation should be minimized in the design. To reduce the separation of the flow, the back sweep angle should be higher. Furthermore, the separation in the rotor with subsonic inlet relative Mach number and normal Reynold's number can be minimized by using a lower limit of the ratio of the outlet to inlet velocity ratio as 0.8.

3. Numerical schemes

In any gas turbine, the flow of the fluid in the compressor is always more unsteady and turbulent, which further makes the design more complicated. For all turbulent flows, the governing equations are the unsteady Navier-Stokes (N-S) equations. But those equations are very much difficult to solve. The following are the governing equations in tensor notation for understanding the basic nature of the equations as presented by [4].

$$\frac{\partial \rho u_i}{\partial t} + \frac{\partial (\rho u_i u_j)}{\partial x_j} = -\frac{\partial (p)}{\partial x_i} + \frac{\partial}{\partial x_j} \left(\mu \left[\frac{\partial (u_i)}{\partial x_j} + \frac{\partial (u_j)}{\partial x_i} \right] \right) - \frac{\partial}{\partial x_i} \left(\frac{2}{3} \mu \frac{\partial (u_k)}{\partial x_k} \right) + S_i \quad (16)$$

Energy equation is given by

$$\frac{\partial \rho \epsilon}{\partial t} + \frac{\partial (\rho u_j \epsilon)}{\partial x_j} = -p \frac{\partial (u_j)}{\partial x_j} + \frac{\partial}{\partial x_j} \left(k \frac{\partial (T)}{\partial x_j} \right) + \emptyset + S_i \quad (17)$$

where

$$\emptyset = \left(\mu \frac{\partial (u_i)}{\partial x_j} \left[\frac{\partial (u_i)}{\partial x_j} + \frac{\partial (u_j)}{\partial x_i} \right] \right) - \left(\frac{2}{3} \mu \left(\frac{\partial (u_k)}{\partial x_k} \right)^2 \right) \quad (18)$$

The main disadvantage in solving N-S equations is computing. This is due to the reason that to represent even very small scale of velocity and pressure fluctuations using N-S equations, temporal resolution is required as the N-S equations are spatial fineness equations. Furthermore, the accuracy and resolution of the scheme reduce due to the increase in accumulative rounding off value errors because of the increase in grid points to achieve fine meshes.

Hence, more accurate schemes are required to achieve the solutions of turbulence models with highest resolution. One such numerical scheme, which can effectively solve the turbulence problems, is Reynolds-averaged Navier-Stokes (RANS) equations.

3.1. Reynolds-averaged Navier-Stokes (RANS) equations

The transport equations should be modified by introducing the components with averaged and fluctuating components for solving the turbulent models. RANS are the equations of the motion of the fluid flow with time-averaged equations. If flow turbulence properties are known, then suitable approximations can be made and high-resolution solutions to the N-S equations can be achieved by solving the RANS equations. The RANS equations in tensor notation are described below [4]

$$\rho \frac{\partial \overline{u_i}}{\partial t} + \rho \frac{\partial (\overline{u_i u_j})}{\partial x_j} = \rho \overline{f_i} + \frac{\partial}{\partial x_j} \left(-\overline{p} \delta_{ij} + 2\mu \overline{S_{ij}} - \rho \overline{u'_i u'_j} \right) \quad (19)$$

where $\overline{S_{ij}} = \frac{\partial (\overline{u_i})}{\partial x_j} + \frac{\partial (\overline{u_j})}{\partial x_i}$ is the mean rate of strain tensor.

The problems of physical engineering, which involve high turbulence when modeled based on RANS equations, are also called statistical turbulence models because the method strongly involves statistical averaging procedure. Consequently, the computational effort to solve RANS equations is very less when compared to the other schemes.

3.2. Governing equations for turbulence models

As discussed in the earlier articles, due to the involvement of more unknown parameters in the modeling of turbulent problems, the accuracy of the model becomes very much difficult and leads to erroneous results. Hence, it is required to develop a numerical procedure to close or converge the system of equations. One such two equations category turbulence model is K-epsilon model [4].

3.2.1. K-epsilon models

In this model, the velocity of the turbulent flow and the length scales are independently calculated using two different equations. The basic model is the standard k-epsilon model [4].

3.2.1.1. Standard k-ε model

The velocity of the turbulence can be calculated by forming a model for the corresponding kinetic energy using the following Eq. (4)

$$\frac{\partial(\rho k)}{\partial t} + \frac{\partial(\rho k u_i)}{\partial x_i} = \frac{\partial}{\partial x_j} \left[\left(\mu + \frac{\mu_t}{\sigma_k} \right) \frac{\partial(k)}{\partial x_j} + P_k + P_b - \rho \epsilon - Y_M + S_k \right] \quad (20)$$

The length scale is represented by ϵ , which is the rate of dissipation and can be calculated from the Eq. (4)

$$\frac{\partial(\rho \epsilon)}{\partial t} + \frac{\partial(\rho \epsilon u_i)}{\partial x_i} = \frac{\partial}{\partial x_j} \left[\left(\mu + \frac{\mu_t}{\sigma_\epsilon} \right) \frac{\partial(\epsilon)}{\partial x_j} + C_{1\epsilon} \frac{\epsilon}{k} (P_k + C_{3\epsilon} P_b) - C_{2\epsilon} \rho \frac{\epsilon^2}{k} + S_\epsilon \right] \quad (21)$$

where P_k represents the turbulence kinetic energy due to the mean velocity gradients, P_b represents the turbulence kinetic energy due to buoyancy, and Y_M represents the contribution of the fluctuating dilatation in compressible turbulence to the overall dissipation rate.

3.3. Mean: Line analysis

The difficulty in the design of the compressor is because it involves two vital phases of the design. One being the 1D design of the compressor and the other is the deep numerical analysis of the design. The difficulty in the first phase is overcome by using mean line theory. In principle, the mean line theory follows the preliminary design that is carried out by neglecting the air flow variations in radial direction and the location of the mean blade radius is considered for the analysis.

3.4. Finite volume discretization method(FVDM)

Any fluid flow problem can be solved with high resolution when the basic concept of continuum is considered during the numerical solution methodology. Finite volume discretization method (FVDM/FVM) is one such methodology developed by McDonald, MacCormack, and Paullay during 1970s. The basic structure of the FVM is as follows:

1. The physical problem is discretized into control volumes.
2. Each control volume is formulated with balance equations involving integrals.
3. Numerical integration is applied to approximate the integral equations.
4. Nodal value interpolations are used to approximate the derivatives and function values.
5. The solution is achieved by assembling the algebraic systems of the discrete control volumes.

4. Case study: Design of the centrifugal compressor

The case study presented in the literature [5] is considered for the present illustration for the design and numerical analysis of the centrifugal compressor. Main objective of the current work is to design a centrifugal compressor capable of delivering a pressure ratio of 5.4 with a mass flow ingestion rate of 5.73 kg/sec. The compressor is targeted to achieve 82% total-total efficiency with design constraints of 280 mm for the impeller and 340 mm for the diffuser in terms of diameter. Also, the compressor has to generate flow parameters relative to functional downstream components.

4.1. Centrifugal compressor specifications

The input parameters of the centrifugal compressor are tabulated in **Table 1**.

S. No	Name of the parameter	Value
1	Mass flow rate	5.73 kg/s
2	Pressure ratio	5.4:1
3	Efficiency	82%
4	Maximum impeller diameter	280 mm
5	Maximum diffuser diameter	340 mm
6	Rotational speed	38,000 rpm
7	Inlet pressure	101,325 Pa
8	Inlet temperature	288.15 K

Table 1. Geometrical and functional specifications of analysed turbo-machinery.

Preliminary design of the impeller and vaned diffuser is carried out by adapting a one-dimensional model approach. Comprehensive design of the centrifugal compressor stage was generated using ANSYS BladeGen module.

4.2. Impeller design

Impeller was designed using guidelines from various sources and in specific from [6]. The number of impeller blade profiles was configured based upon a choice iterative method aiming for passage flow devoid of heavy separations, and it was streamlined and set at 19 blade profiles. A rotational speed of 38,000 rpm was set for the impeller.

4.2.1. Impeller inlet

Hub-to-tip ratio is a key owing to the secondary losses, which occurs in the flow regions near the end walls. The presence of any undesirable circulatory or cross-flows develops on account of rapid and steep flow turning through the blade channel accounting for annulus wall boundary layers. Therefore, impeller hub-to-tip ratio was set at 0.35 within the range of 0.3–0.6 prescribed in [6].

Mach number at impeller inlet was set at 0.65. Impeller inlet blade angles were setup by PCA Engineer's Vista CCD tool analytical calculations. The inducer leading edge angles are 35, 56, and 63°, respectively, with incidences of 11.8, 3.7, and 0.1° at hub, mean, and shroud, respectively. Leading edge of the impeller was defined by an elliptic ratio of 6.

4.2.2. Impeller outlet

Impeller backswept angle was set at 0° to minimize impeller diameter, and a lean angle of 30° was also incorporated into the design. Impeller exducer height at 13.7032 mm and impeller diameter of 257.954 mm were set by PCA Engineer's Vista CCD software impeller trailing edge that was defined as a square cutoff.

Data generated by Vista CCD tool were used to generate a 3D computer-aided design model of impeller. Inlet portion of 50 mm and horizontal was designed satisfying the diameter constraint of 280 mm.

4.3. Design of vaned diffuser

In order to obtain higher pressure ratio in a radial diffuser, the diffusion process has to be achieved across a relatively shorter radial distance. This requires the application of vanes, which provide greater guidance to flow inside diffusing passages. The vaned diffuser was designed by observing various flow parameters reflected at impeller exit after performing numerical simulations.

To circumvent flow separation, divergence of diffuser blade passages in vaned diffuser ring can be kept small by incorporating a large number of vanes. However, this can lead to higher friction losses. Thus, an optimum number of diffuser vanes must be employed and ensure flow passage divergence not to exceed 12°. Thus, final diffuser design contains 30 blades. The

diffuser vane leading edge was at a radius of 136 mm, and the trailing edge of the diffuser vanes was set at a radius of 166 mm. Blade inlets and outlet angle were set at 64° . The leading and trailing edge were defined by an elliptic ratio of 6 and a radius of 0.25 mm.

Another method to prevent very steep velocity gradients at diffuser entry is by providing a small vane-less space ($0.05d_2$ – $0.1d_2$) between impeller exit and diffuser entry. Therefore, a vane-less space of 7.023 mm was allotted.

The design methodology adapted was mainly focused aiming at decrease in Mach number and flow angle at diffuser exit by satisfying a diameter constraint of 340 mm.

4.4. Meshing strategy

Adapted meshing strategy for both the impeller and diffuser fluid domains was achieved using Ansys-Turbo grid module. A total node count of $7e-01$ million was setup for the CFD solver, as it allows generation of refined quality hexahedral meshes required for the blade passages in turbo-machinery.

4.5. Numerical setup

In the present context, ANSYS tool is setup with pressure-based solver simulating a steady, three-dimensional viscous flow fields using complete set of Navier-Stokes code solving for Reynolds-averaged Navier-Stokes equations based on finite volume discretization method [7]. A high-resolution scheme is used to solve for continuity, momentum, energy, and state equations implementing a standard k - ϵ turbulence model. Individual compressor stage characteristics were generated by performing simulations varying back pressures. Firstly, near choke condition, flow points are run to reduce static back pressure values and later, solutions are restarted with incrementally increasing the static back pressure to compute intermediate points on constant speed line running toward stall.

Inlet boundary conditions: At compressor inlet, a constant total pressure and total temperature conditions are imposed with a turbulence intensity of 1% and flow direction is marked normal to the inlet plane.

Outlet boundary conditions: At the outlet, an average static pressure boundary condition is implemented. Also, a circumferential symmetry condition is imposed on corresponding periodic surfaces with air-fluid medium setup as an ideal gas. A counter-rotating wall boundary is given at the impeller shroud.

5. Results and discussions

5.1. Flow analysis: Velocity triangles

Entry and exit velocity triangles for impeller blades along the radial section with backward swept blades ($\beta_2 = 0^\circ$) are plotted as shown in **Figure 2(a)** and **(b)**. Inlet Mach number is

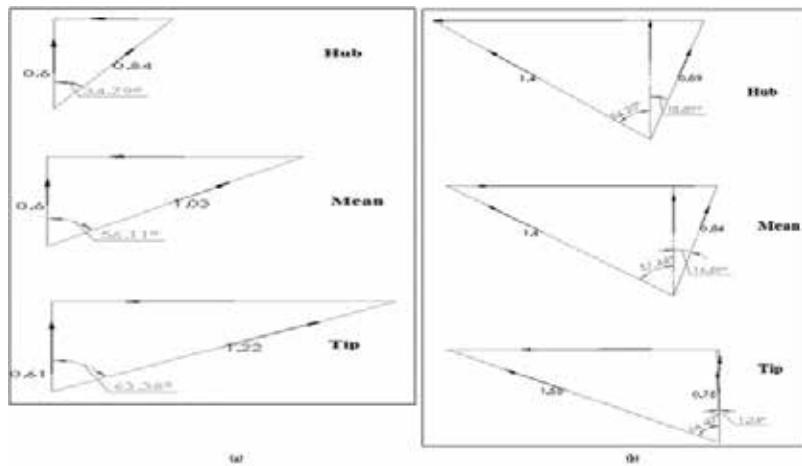


Figure 2. (a) Inlet Mach triangles for impeller blades only in radial section with backward swept blades, $\beta_2 = 0^\circ$ with zero swirl at entry; (b) exit Mach triangles for impeller blades only in radial section with backward swept blades, $\beta_2 = 0^\circ$ with zero swirl at entry.

within the recommended range of 0.4 to 0.6. A maximum flow angle of 64° at the impeller exit and a Mach number of 1.4 are observed.

5.2. Blade to blade contours

Blade to blade contours for absolute Mach number at hub, mean and tip sections are presented in below **Figure 3 (a)** and **(c)**.

The absolute Mach number plots show no supersonic regions of significant volume. The hub section plot show supersonic flow near suction surface of the blade and the tip Section span plots show a supersonic region on the pressure surface of the blade. The presence of a mild bow shock in front of the leading edge of the blade can also be observed. The Mach number is around 0.6 at the leading edge of the impeller. The following plots shown in **Figure 4(a–c)** are the relative Mach number contours for the hub, mean and tip.

The Mach number at the leading edge in the relative frame is near 1.2 for the leading edge of the impeller blade. The regions of low Mach number can also be observed, although it is in the relative frame.

5.3. Circumferential contour

Circumferential relative Mach number contours at the impeller inlet and outlet are presented in **Figure 5(a)** and **(b)**. The higher color temperature indicates higher values of relative Mach number.

From above **Figure 5(a)**, corresponding inlet relative Mach number contours shows that passage area is lying in the supersonic regime. The **Figure 5(b)** shows the effect of lean angle of impeller blades.

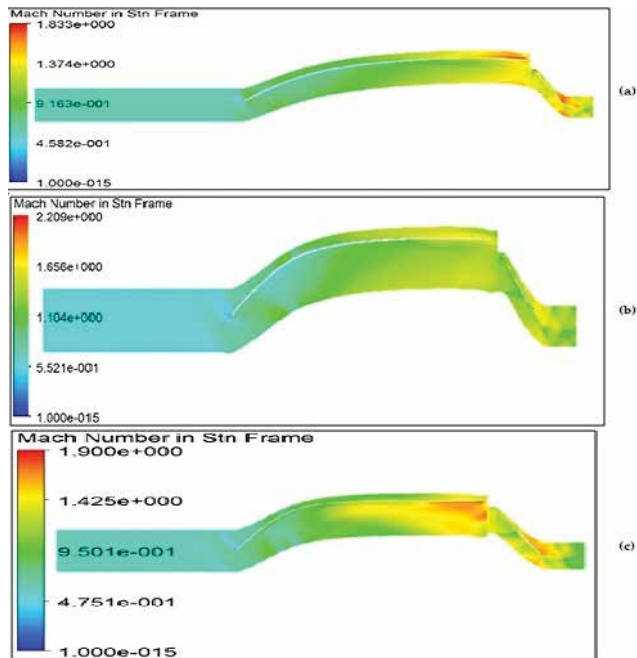


Figure 3. (a) Hub section absolute Mach number contour, (b) mean section absolute Mach number contour, and (c) tip section absolute Mach number contour.

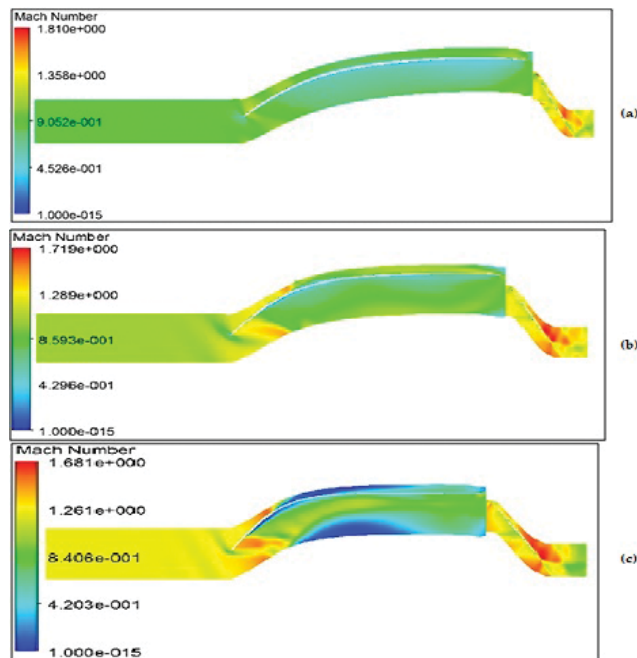


Figure 4. (a) Hub section relative Mach number contour, (b) mean section relative Mach number contour, and (c) tip section relative Mach number contour.

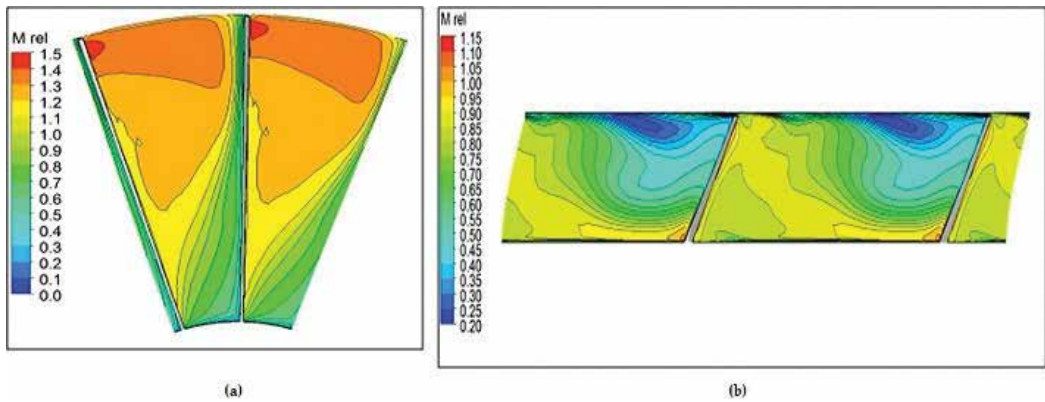


Figure 5. (a) Inlet relative Mach contour and (b) outlet relative Mach contour.

5.4. Performance characteristics

Individual component as well as stage performance as a whole can be gauged by different methods.

5.4.1. Impeller performance characteristics

Impeller performance was evaluated by Vista-CCM and their corresponding plots are shown in below **Figure 6** represents the variation of pressure ratio with varying mass flow ingested through compressor.

Figure 7 represents the variation of polytropic efficiency with varying mass flow ingested to quantify differential pressure changes occurring through compressor stage.

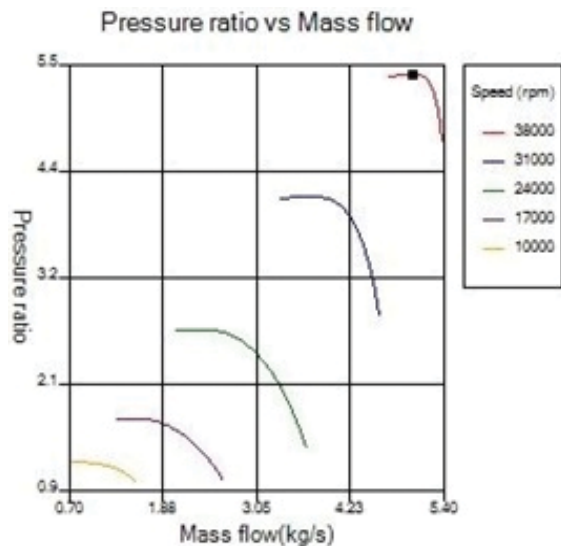


Figure 6. Mass flow versus pressure ratio by Vista CCM.

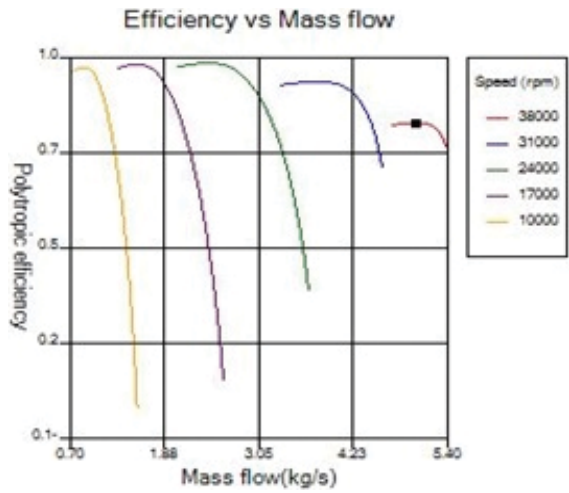


Figure 7. Mass flow versus polytropic efficiency by Vista CCM.

5.4.2. Stage performance characteristics

Compressor stage performance can be evaluated by obtaining individual characteristic curves as plotted in Figures 8–10. The compressor performance is presented only for a functional design speed of 38,000 rpm. Variables like total pressure ratio, adiabatic efficiency, and power requirement of compressor are plotted against varying mass flow rate as shown in Figures 8–10. The maximum achieved pressure ratio is 5.4 while allowing a fixed mass flow rate of 5.73 kg/s.

A peak efficiency of 75.8% is obtained at a pressure ratio of 5.73. The design point efficiency is predicted to be around 74.5%. Possible reasons for loss in efficiency may be owing to primary and secondary losses.

The impeller works with power a power requirement in the range of 1350–1360 kW on its 100% speed line over the range of total pressure ratio 4.5–5.4.

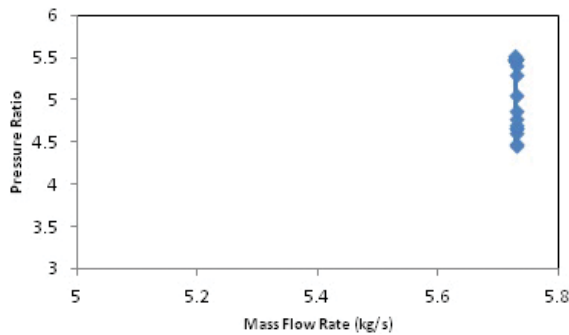


Figure 8. Pressure ratio versus mass flow.

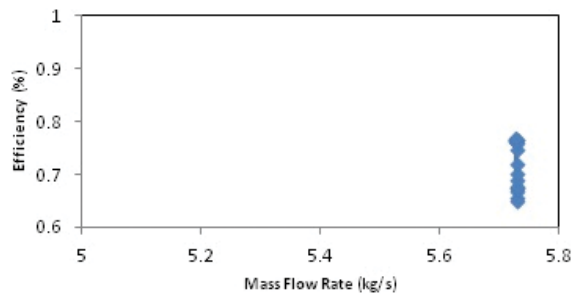


Figure 9. Efficiency versus mass flow.

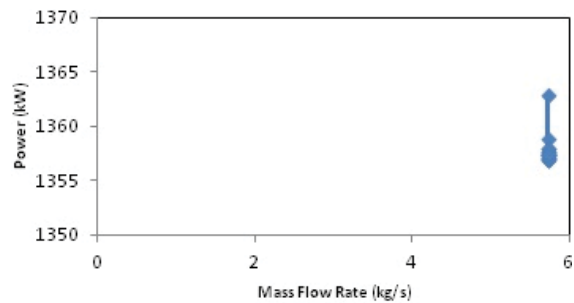


Figure 10. Power versus mass flow.

5.5. Validation studies

The obtained results presented in the Figures 11 and 12 shows state-of-the-art validity of current design with design performance work stated in the literature.

The design methodology adapted in current framework uses a backswept angle of 0° and does not follow the design trend established for validation in above figures. Above deviations in the performance trends can be attributed to the structural size and rotational specifications of various stages.

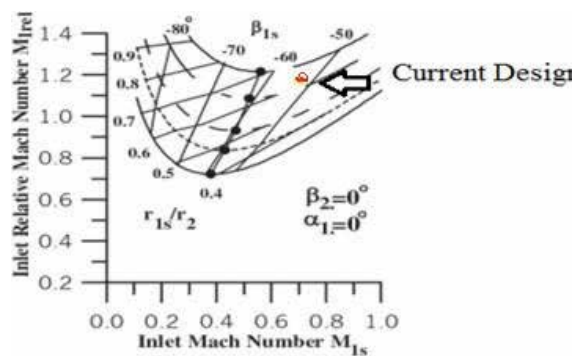


Figure 11. Comparison with compressor in [8].

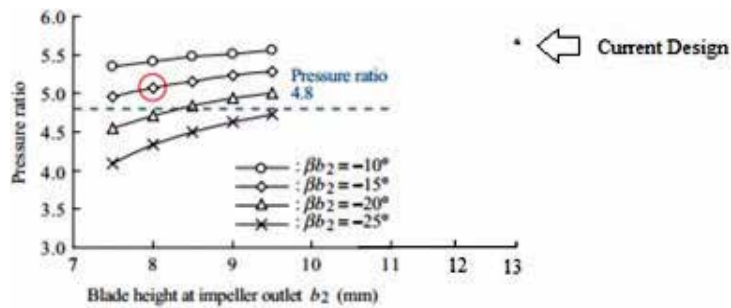


Figure 12. Comparison with trend established in paper [9].

6. Conclusion

The aero-thermodynamic design of centrifugal compressor turbo-machinery for a small gas turbine engine has been successfully carried out and sequentially investigated for its performance at design speed. The 3D computational analysis infers that the designed centrifugal compressor generates a pressure ratio of 5.4 with an ingested mass flow rate of 5.73 kg/s at a fixed rotational speed of 38,000 rpm. The targeted flow parameters are in well agreement with centrifugal turbo-machinery within the diameter specification of 280 and 340 mm for impeller and diffuser, respectively. The total-to-total efficiency of designed centrifugal compressor is 81% against a target of 82%. The swirl at impeller is reduced from 77 to 27.5° by the diffuser, thus reducing the primary losses viz., flow separation and thereby has the capability of experiencing lesser stall angles.

Author details

Chellapilla V K N S N Moorthy^{1*} and Vadapalli Srinivas²

*Address all correspondence to: krishna.turbo@gmail.com

1 Institute of Aeronautical Engineering, Hyderabad, Telangana, India

2 GITAM University, Vishakhapatnam, Andhra Pradesh, India

References

- [1] Rohlik HE. Radial inflow turbines (in turbine design and application). NASA Special Publication SP 290. 1975;3:31-58
- [2] Rodgers C. Specific speed and efficiency of centrifugal impellers. In: Performance Prediction of Centrifugal Pumps and Compressors. New York: ASME; 1980. pp. 592-603

- [3] Wilson DG, Korakianitis TH. The Design of High Efficiency Turbomachinery and Gas Turbines. 2nd ed. New Jersey, United States: Prentice Hall; 2014
- [4] Srinivasa Rao P: Modeling of turbulent flows and boundary layer. In: Hyoung Woo Oh, editor. Computational Fluid Dynamics. London, United Kingdom: InTech; 2010. p. 285-306. ISBN: 978-953-7619-59-6
- [5] Moorthy CHVKNSN, Bharadwajan K, Srinivas V. Computational studies on aerothermodynamic design and performance of centrifugal turbo-machinery. International Journal of Mechanical Engineering and Technology. 2017;**8(5)**:320-333
- [6] Philip P Walsh, Fletcher P. Gas Turbine Performance. 2nd ed. Oxford: John Wiley & Sons; 1998. p. 178-186
- [7] Moorthy CHVKNSN, Srinivas V, Prasad VVSH. Computational analysis of a CD nozzle with 'SED' for a rocket air ejector in space applications. International Journal of Mechanical and Production Engineering Research and Development. 2017;**7(1)**:53-60
- [8] Zahed AH, Bayomi NN. Design procedure of centrifugal compressors. ISESCO Journal of Science and Technology. 2014;**10(17)**:77-91
- [9] Hideaki T, Masaru U, Hirata KT. Aerodynamic design of centrifugal compressor for AT14 turbocharger. IHI Engineering Review. 2010;**43(2)**:70-76

Modeling and Simulation of Problems of Fundamental Physics

Susceptibility of Group-IV and III-V Semiconductor-Based Electronics to Atmospheric Neutrons Explored by Geant4 Numerical Simulations

Daniela Munteanu and Jean-Luc Autran

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.71528>

Abstract

New semiconductor materials are envisaged in numerous high-performance applications for which the expected device or circuit performances cannot be achieved with silicon. In this context of growing use of new and specific semiconductors, the question of their susceptibility to natural radiation, primarily to atmospheric neutrons, is posed for high-reliability-level application domains. This numerical simulation work precisely examines nuclear events resulting from the interaction of atmospheric neutrons at the terrestrial level with a target layer composed of various group-IV and III-V semiconductor materials including silicon, germanium, silicon carbide, carbon-diamond, gallium arsenide, and gallium nitride materials. Using extensive Geant4 simulations and in-depth data analysis, this study provides an accurate and fine comparison between the neutron interaction responses of these different semiconductors in terms of nuclear processes, recoil products, secondary ion production, and fragment energy distributions. Implications of these results on the rate of single-event transient effects at the device or circuit level are also discussed.

Keywords: terrestrial cosmic rays, atmospheric neutrons, neutron-semiconductor interactions, group-IV semiconductors, III-V semiconductors, nuclear reactions, Geant4, numerical simulations, radiation effects on electronics, single-event effects

1. Introduction

Modern electronics and optoelectronics are increasingly based on specific semiconductor materials other than silicon, including some elemental or compound semiconductors like germanium, silicon-germanium, and silicon carbide (group-IV) and III-V alloys [1]. These new materials are

attracting strong interest because of their better electronic transport, optical, or high-frequency properties than Si; they can be envisaged in numerous high-performance applications (e.g., “More than Moore” microelectronics and beyond CMOS, extreme environments, high temperatures, or high-speed electronics) for which the expected device or circuit performances cannot be achieved with silicon. In such a context of growing use of new and specific semiconductors, the question of their susceptibility to natural radiation, primarily to atmospheric neutrons at ground level, is posed for high-reliability-level application domains. A special attention should be particularly paid to low-bandgap materials (Ge and most of III-V materials), envisaged as channel replacement for MOSFETs and steep switching tunnel FETs for low voltage application [2], due to their low ionization energy susceptible to amplify charge generation from sea-level neutron radiation.

Following a methodology previously developed for the study of neutron-silicon interactions [3], the present work precisely examines nuclear events resulting from the interaction of atmospheric neutrons at the terrestrial level with a target layer composed of Si, Ge, SiC, C (diamond), GaAs, and GaN materials and representative of the whole sensitive volume of a typical integrated circuit. To perform this task, we constructed using the Geant4 toolkit [4, 5] a specific source of atmospheric neutrons and compiled large databases of neutron-semiconductor interaction events corresponding to tens of thousands of nuclear reactions. Details of these simulations and database compilation are given in Section 2. Section 3 presents a detailed analysis of obtained databases in terms of nuclear processes, recoil products, secondary ion production, and fragment energy distributions. Finally, Section 4 discusses the implications of these results on the rate of single-event transient effects at the device or circuit level.

2. Simulation details

2.1. Semiconductor materials

For the purpose of this study, we considered different group-IV and III-V semiconductor materials including silicon, germanium, silicon carbide, carbon-diamond, gallium arsenide, and gallium nitride materials. **Table 1** summarizes the main physical and atomic properties of these materials [6] notably in terms of bandgap value, bulk density, average energy for creation of an electron-hole pair and energy threshold to deposit 1.8 fC in the considered material. This last quantity corresponds to a silicon recoil nucleus in silicon material with a kinetic energy of 40 keV, considered in [3] as the lowest deposited charge susceptible to induce a soft error in a bulk 65 nm SRAM memory. For standardization purpose and comparison with results presented in [2], we adopted here the same lowest deposited charge (1.8 fC corresponds to 11,250 electrons) that leads to $E_{th} = 11,250 \times E_{eh}$. In other words, all neutron-induced products with energies below E_{th} will not be considered in the databases (see paragraph 2.3).

2.2. Atmospheric source of neutrons

The different semiconductor materials given in **Table 1** have been irradiated with atmospheric neutrons in Geant4 simulations described below (paragraph 2.3) and schematically represented in the inset of **Figure 1**. For memory, the interaction of primary cosmic rays with the Earth’s top

Properties (300 K)	Si	Ge	C (diamond)	4H-SiC	GaN	GaAs
Semiconductor group	IV			III-V		
Atomic number	14	32	6	14/6	31/7	31/33
Bandgap (eV)	1.124	0.661	5.47	3.23	3.39	1.424
Structure	Diamond (cubic)			Wurtzite (hexagonal)		Zinc Blende (cubic)
Lattice constant (Å)	5.43	5.65	3.567	3.73 10.053	3.186 5.186	5.653
Density (g/cm ³)	2.329	5.3267	3.515	3.21	6.10	5.32
Atoms (/cm ³)	5.0×10^{22}	4.42×10^{22}	1.76×10^{23}	5.0×10^{22}	8.7×10^{22}	4.5×10^{22}
Average energy (E_{eh}) for creation of an electron-hole pair (eV)	3.6	2.9	12	7.8	8.9	4.8
Energy threshold (E_{th}) to deposit 1.8 fC (keV)	40	33	135	58	100	47

Table 1. Main structural, atomic, and electronic properties of the different group-IV and III-V semiconductor materials considered in the present study. (Data partially from Ref. [6]).

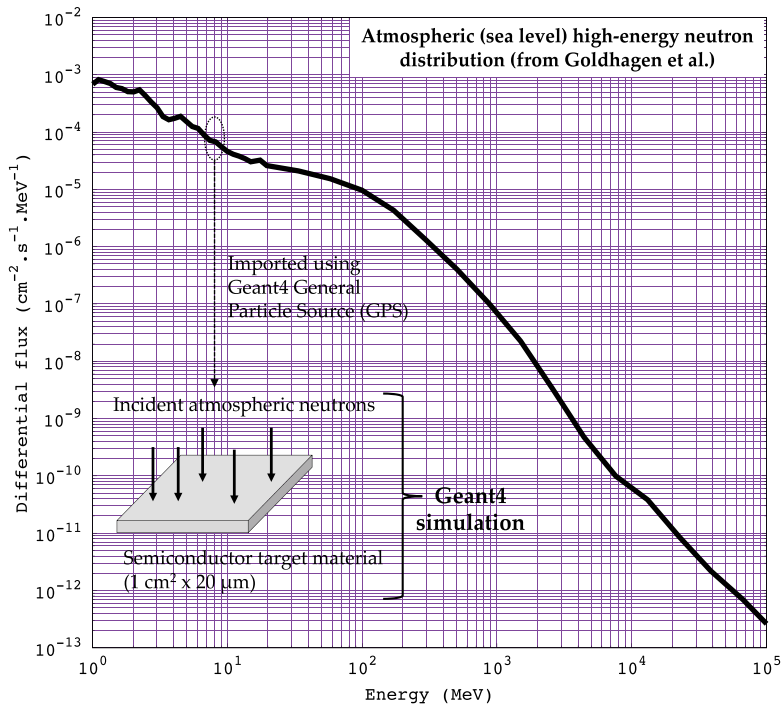


Figure 1. Differential flux of cosmic-ray induced high-energy neutrons as measured by Gordon and Goldhagen et al. using a multielement Bonner sphere spectrometer on the roof of the IBM T. J. Watson Research Center in Yorktown Heights, NY [6]. *Inset:* Schematic representation of the neutron irradiation simulated using Geant4 in this work.

atmosphere is at the origin of atmospheric showers that produce secondary particles down to the sea level. After muons, the next most abundant secondary particles at the sea level are neutrons. High-energy neutrons (typically above 1 MeV) represent by far the main threat to electronics at the ground level because these particles being not charged are very invasive and can penetrate deeply in circuit materials where they can interact with atoms to produce charged products (recoil nuclei or secondary ions).

To emulate the atmospheric neutron source, the differential flux of cosmic-ray induced high-energy neutrons measured by Gordon and Goldhagen et al. in Yorktown Heights [7] has been considered as the reference input spectrum [8]. This distribution, shown in **Figure 1**, was imported in the Geant4 general particle source (GPS) library [9] to randomly generate incident neutrons mimicking the natural sea-level neutron background.

2.3. Geant4 options, models, and simulation runs

The neutron interaction databases for the different semiconductor materials listed in **Table 1** have been computed in this work using Geant4 version 4.9.4 patch 01. The list of physical processes employed in simulation is based on the standard package of physics lists QGSP_BIC_HP [10]. Concerning the hadronic interactions, in QGSP group of physics lists, the quark gluon string model is applied for high-energy (above ~12 GeV) interactions of protons, neutrons, pions, kaons, and nuclei. The high-energy interaction creates an excited nucleus, which is passed to the precompound model describing the nuclear de-excitation. Nuclear capture of negative particles is simulated within the chiral invariant phase space (CHIPS) model. QGSP_BIC_HP list includes binary cascade for primary protons and neutrons with energies below ~10 GeV and also uses binary light ion cascade for inelastic interaction of ions up to few GeV/nucleons with matter. The complete list of the Geant4 classes that we considered for our neutron simulations is summarized in **Table 2**.

For each semiconductor target, a simulation run consists in the generation of 100 millions (10^8) of primary incident neutrons on a 20- μm -thick material layer perpendicular to its surface (1 cm^2).

Neutron process	Energy	Geant4 model	Dataset
<i>Elastic</i>	< 20 MeV	G4NeutronHPElastic	G4NeutronHPElasticData
	> 20 MeV	G4LElastic	-
<i>Inelastic</i>	< 20 MeV	G4NeutronHPInelastic	G4NeutronHPInelasticData
	[20 MeV, 10 GeV]	G4BinaryCascade	-
	[10 GeV, 25 GeV]	G4LENeutronInelastic QGSP	-
	[12 GeV, 100 TeV]		-
<i>Fission</i>	< 20 MeV	G4NeutronHPFission G4LFission	G4NeutronHPFissionData
	> 20 MeV		-
<i>Capture</i>	< 20 MeV	G4NeutronHPCapture G4LCapture	G4NeutronHPCaptureData -
	> 20 MeV		

Table 2. List of the different Geant4 classes considered in the present simulation flow for the description of neutron-matter interactions.

All particles characterized by tiny (i.e., insignificant) ionizing properties and thus negligible impact on electronics in terms of electron-pair generation and single events have been excluded from the simulation. Consequently, the computed databases exclude electrons, positrons, gamma photons, pions, and mesons and only contain protons, alpha particles, and ionizing products with $Z > 2$.

The target dimensions have been chosen to match the typical dimensions of the sensitive volume of an integrated circuit. In particular, the thickness has been fixed to 20 μm because the electrical charge generated by a reaction product beyond 20 μm would not drift or diffuse to the active area (i.e., the sensitive region of the circuit located close to the semiconductor surface) and, consequently, would not play any role in the occurrence of single events.

3. Simulation results

3.1. Main characteristics of the databases

An example of the computed databases is given in **Table 3** for a target composed of GaAs material and subjected to atmospheric neutron irradiation. As illustrated in this sample composed of nine events, each neutron-interaction event is described in two or several successive lines: the first line gives the event number, the energy of the incident neutron at the origin of the event, the Cartesian coordinates of the reaction vertex, and the number of secondary products generated; the following lines give for each secondary product the nature of this product, its mass and atomic numbers, its initial energy just after emission/production, and the Cartesian coordinates of its momentum. In the first line of the database sample shown in **Table 3**, the event #11369 corresponds to a neutron elastic interaction with an As nucleus of the semiconductor lattice, giving a recoil As with an energy of 58 keV.

Table 4 summarizes the main size characteristics of the computed databases for the six semiconductor materials listed in **Table 1** and subjected to 100 million of atmospheric neutrons each. Five size parameters (or metrics) are reported in **Table 4** to quantify the different databases in terms of neutron-semiconductor interactions: the number (fraction) of elastic and inelastic events, the total number of events, the total number of generated secondary products with energy above the energy threshold E_{th} given in **Table 1**, and, finally, the average number of secondary ions produced in the case of inelastic events. **Figures 2** and **3** also graphically represent the fractions of elastic and inelastic events and the total number of interaction events and generated secondary products, respectively. All these results lead to the following remarks:

- Silicon exhibits the lowest total number of interaction events, and, except the particular case of carbon-based materials (diamond and SiC), GaN shows the highest event rate with more than 50% of the supplementary events with respect to Si.
- This result can be explained by the fact that neutron nuclear interaction probability with the elements is even higher than the atomic number Z is. This means that the susceptibility of the materials is all the greater as their atomic number is high [11].

- Diamond shows a very different behavior than the other materials since it is an excellent neutron moderator (better than graphite). This explains the extremely elevated number of elastic events as compared to other materials: $\times 8$ with respect to Ge, $\times 7$ with respect to GaAs, and $\times 4$ with respect to Si.
- Silicon carbide, which can be viewed as a “mixture” of Si and C at the atomic level, shows an intermediate behavior between Si and C with a quasi $\times 2$ number of elastic events due to the presence of C.

Figure 2 shows the fraction of elastic and inelastic events for the different semiconductor materials. We can distinguish three different behaviors with low ($< 30\%$, Ge and GaAs), intermediate

Event Number	Energy (MeV) of the incident neutron	(x, y, z) coordinates of the vertex of the interaction event	Number of products
11369	1 1.938009e+00	-3.123114e+00 2.128556e+00 -2.043811e-03	1
As75[0.0]	75 33 5.865383e-02	4.466775e-01 4.319914e-01 7.834939e-01	
Product	A Z Energy (MeV)	(ux, uy, uz) coordinates of the momentum vector	
11370	1 2.383369e+02	-1.712552e+00 -2.155406e+00 -2.719637e-03	3
proton	1 1 8.546443e+01	-5.270216e-01 -1.286384e-01 8.400597e-01	
alpha	4 2 1.383653e+01	2.459101e-02 -1.145710e-01 -9.931107e-01	
Ni62[0.0]	62 28 2.160926e+00	8.943092e-02 -1.490662e-01 9.847748e-01	
11371	1 2.638142e+01	-3.206409e+00 -4.575562e+00 9.494563e-03	1
Ga68[0.0]	68 31 6.260841e-01	1.892272e-01 -2.680785e-01 9.446306e-01	
11372	1 2.074821e+02	4.896602e+00 -4.752904e+00 6.255025e-04	2
proton	1 1 1.893323e+02	1.776401e-01 2.641910e-02 9.837408e-01	
Zn68[0.0]	68 30 1.562386e-01	-9.659646e-01 9.698920e-03 2.584925e-01	
11373	1 7.875999e+00	-1.735542e+00 4.719762e+00 -5.800941e-03	1
Ga69[0.0]	69 31 2.959972e-01	3.094430e-02 9.454495e-02 9.950395e-01	
11374	1 1.114444e+02	4.558722e+00 -4.191099e+00 9.470209e-03	2
proton	1 1 6.645712e+00	-4.477549e-01 -6.883188e-01 -5.707301e-01	
Zn64[0.0]	64 30 3.507131e+00	3.622407e-01 1.244171e-01 9.237435e-01	
11375	1 1.541558e+02	2.595860e-01 1.292310e+00 5.783486e-03	2
proton	1 1 3.268971e+01	-3.472904e-02 -3.481138e-01 9.368088e-01	
Ge74[0.0]	74 32 4.819644e-01	-9.020135e-01 3.033525e-01 -3.071628e-01	
11376	1 3.974678e+01	6.775934e-01 4.367106e+00 -9.091089e-03	2
alpha	4 2 7.800831e+00	7.650253e-01 3.799136e-01 5.200019e-01	
Cu64[0.0]	64 29 1.066644e+00	-5.057675e-01 -4.010571e-01 7.637751e-01	
11377	1 2.877963e+02	-1.109289e+00 1.109673e+00 -8.077168e-03	4
proton	1 1 6.803890e+01	2.508779e-01 -6.180162e-01 7.450612e-01	
proton	1 1 1.804781e+01	-3.496855e-01 -7.349503e-01 -5.810061e-01	
deuteron	2 1 1.321282e+01	-5.108815e-01 -4.381425e-01 7.396156e-01	
Zn66[0.0]	66 30 9.454173e-01	9.760912e-01 2.947766e-04 2.173610e-01	

Table 3. Sample (nine events) extracted from the computed database corresponding to GaAs target material subjected to atmospheric neutron irradiation.

Target material	Si	Ge	C (diamond)	4H-SiC	GaN	GaAs
Number of elastic events (fraction, %)	7918 (59.1)	3788 (27.7)	30,536 (79.2)	13,771 (65.1)	8303 (39.1)	4119 (28)
Number of inelastic events (fraction, %)	5482 (40.9)	9898 (72.3)	7981 (20.8)	7390 (34.9)	12,955 (60.9)	10,603 (72)
Total number of events	13,400	13,686	38,517	21,161	21,258	14,722
Total number of secondary products (E > E_{th})	18,989	20,527	48,357	29,067	32,011	21,758
Average number of secondary ions produced in inelastic events	2.02	1.69	2.23	2.07	1.83	1.66

Table 4. Main characteristics of the neutron event databases generated for the six semiconductor materials considered in this work.

(40–60%, Si and GaN), and high (> 60%, C and SiC) elastic event rates; the presence of low-Z elements such as C and N, respectively, in SiC and GaN leads to increase elastic interactions in these last materials with respect to the elastic rates observed for Si and GaAs.

In addition, **Figure 3** shows the total number of interaction events (elastic + inelastic) and the total number of secondary products generated in the different targets. For Si, Ge, and GaAs, this last quantity is in the same order of magnitude (with Si < Ge < GaAs that follows the rule that neutron nuclear interaction probability with the elements is even higher than the atomic number is). For GaN and SiC and as previously mentioned, the presence of low-Z elements, C and N, respectively, increases the number of elastic events and indirectly increases the total number of secondary products. Finally, carbon-diamond material dominates this comparison in terms of the number of events/products due to its high power of neutron moderation.

3.2. Nature of the secondary products

Figure 4 shows the number of events as a function of the number of secondary products generated during the different interaction events. We call this last quantity secondary product

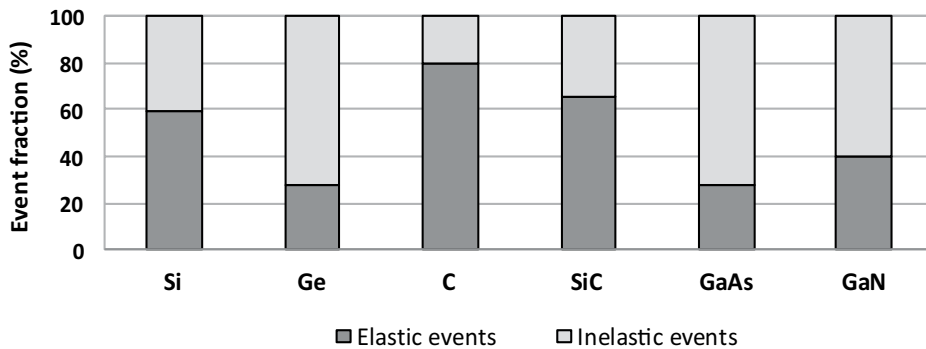


Figure 2. Fraction of elastic and inelastic events for the six semiconductor targets irradiated with atmospheric neutrons.

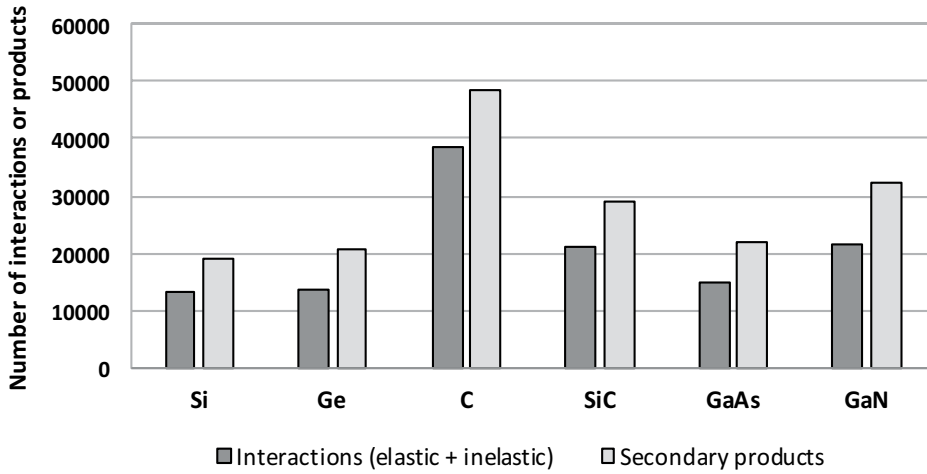


Figure 3. Total numbers of interaction events and secondary products referenced in the different computed databases.

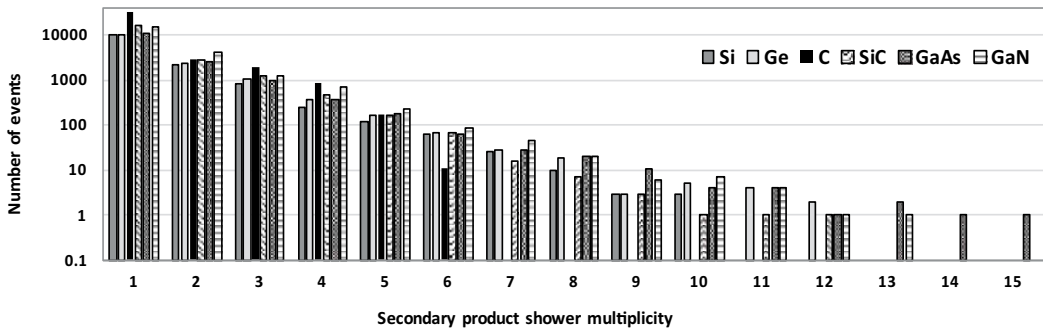


Figure 4. Number of events as a function of the number of secondary products (also called secondary product shower multiplicity) for the six semiconductor materials (for 100 million atmospheric incident neutrons on a volume target of 1cm² x 20 microns).

shower multiplicity (M) because each event can be considered at the origin of a shower of ionizing products: a multiplicity of one (i.e., one product emitted) corresponds to elastic events; larger multiplicities correspond to the production of two or more secondary products during nuclear reactions.

From **Figure 4**, we can formulate the following observations:

- For all targets, the number of reactions monotonously decreases when increasing M . This number of reactions goes to zero above $M = 6$ for carbon (diamond), $M = 10$ for silicon, $M = 12$ for germanium and silicon carbide, $M = 13$ for gallium nitride, and $M = 15$ for gallium arsenide.
- For $M > 9$, the event statistic is therefore relatively weak and may be dependent of the number of incident neutrons. Pushing the statistics beyond 100 million neutrons may give slightly different results for these limits in terms of secondary product shower multiplicity.

- Multi-fragment reactions with $M \geq 4$ represent a small but non-negligible part of the events for all semiconductors except for diamond: 3.6% for Si, 4.9% for Ge, 4.6% for GaAs, 5.0% for GaN, and 3.2% for SiC against only 2.8% for C. These large multiplicity events are important because they can produce a single event that potentially impacts several sensitive volumes in a component or a circuit. With respect to silicon, Ge, GaAs, and GaN show a slightly higher probability for such large multiplicity events.

Figures 5 and 6 provide a detailed analysis of the production of these secondary ions, target per target, and as a function of the atomic number of the products. Z is ranging from 1 (proton) to the highest atomic number of the element(s) present in the target material: 6 for carbon, 14 for silicon, 32 for germanium, 31 for GaN, and 33 for GaAs.

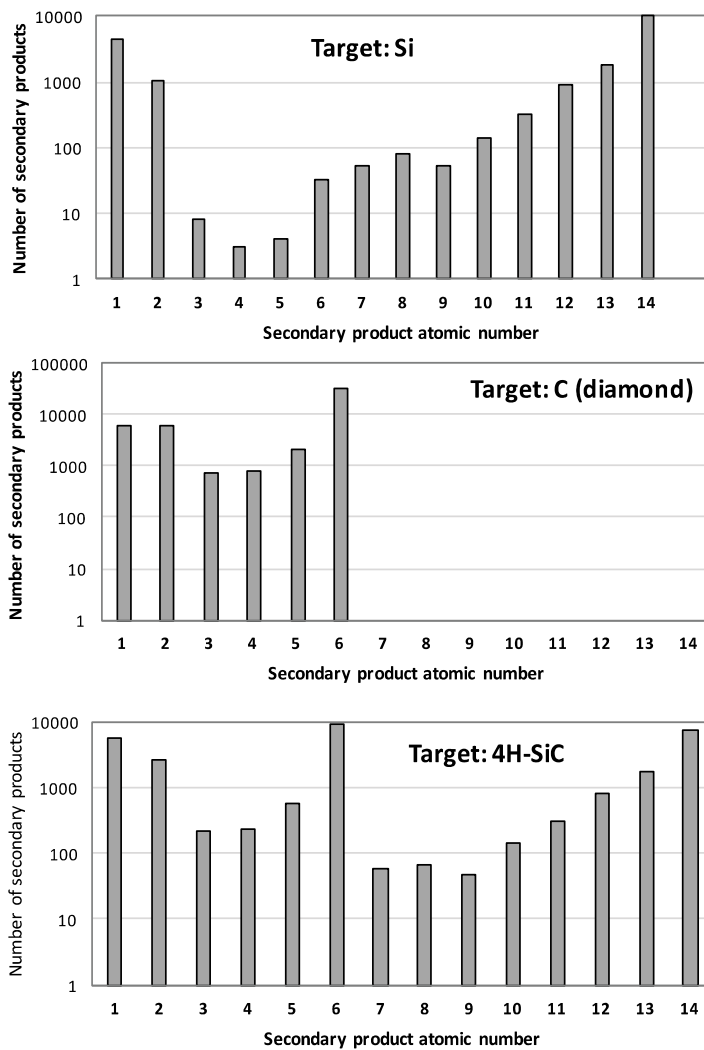


Figure 5. Number of secondary products produced in Si, C (diamond), and SiC as a function of their atomic number (for 100 million atmospheric incident neutrons on a volume target of $1 \text{ cm}^2 \times 20 \text{ microns}$).

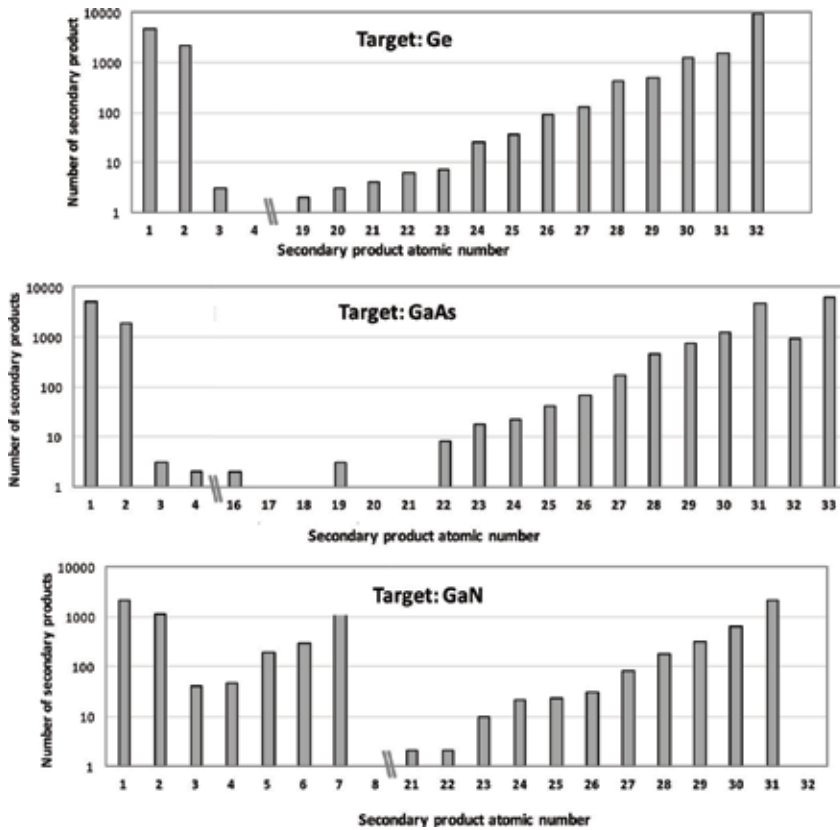


Figure 6. Number of secondary products produced in Ge, GaAs, and GaN as a function of their atomic number (for 100 million atmospheric incident neutrons on a volume target of $1 \text{ cm}^2 \times 20 \text{ microns}$).

For the six materials, the most frequent produced secondary particles are the recoil products due to neutron elastic interactions with the nuclei of the semiconductor lattice, followed by protons in the second position and alpha particles in the third one (in the case of compound materials, recoil nuclei concern the two different species: Si and C recoil products for SiC, Ga and As for GaAs, and Ga and N for GaN). All the other products are systematically less produced than these three categories of products.

This observation justifies why in the following analysis (paragraph 3.3), all produced particles will be divided in four classes: recoil products, protons, alpha particles, and other products.

Also shown in **Figure 5**, beryllium ($Z = 4$) is the less produced product for silicon target, lithium ($Z = 3$) for diamond, and fluor ($Z = 9$) for SiC. Finally, note that in **Figure 6**, an axis-break has been introduced because no secondary product is created between $Z = 3$ and 19 for Ge, between $Z = 4$ and 22 (except $Z = 16$ and 19) for GaAs, and between $Z = 7$ and 21 for GaN. For these large Z elements (31 for Ga, 32 for Ge, and 33 for As), the absence of fragments in these ranges of Z

comes from the limited number of high-energy incident neutrons, due to the $1/E$ nature of the spectrum of **Figure 1** and to the finite value of the high-energy limit (10^5 MeV) considered for simulation; certain reaction channels being not present in the computed databases.

3.3. Detailed analysis in terms of energy, LET, and range

Figures 7–10 provide a detailed analysis of the secondary ions produced during neutron interactions in the different target materials in terms of initial energy (when products are released), linear energy transfer (LET), and range in the target material.

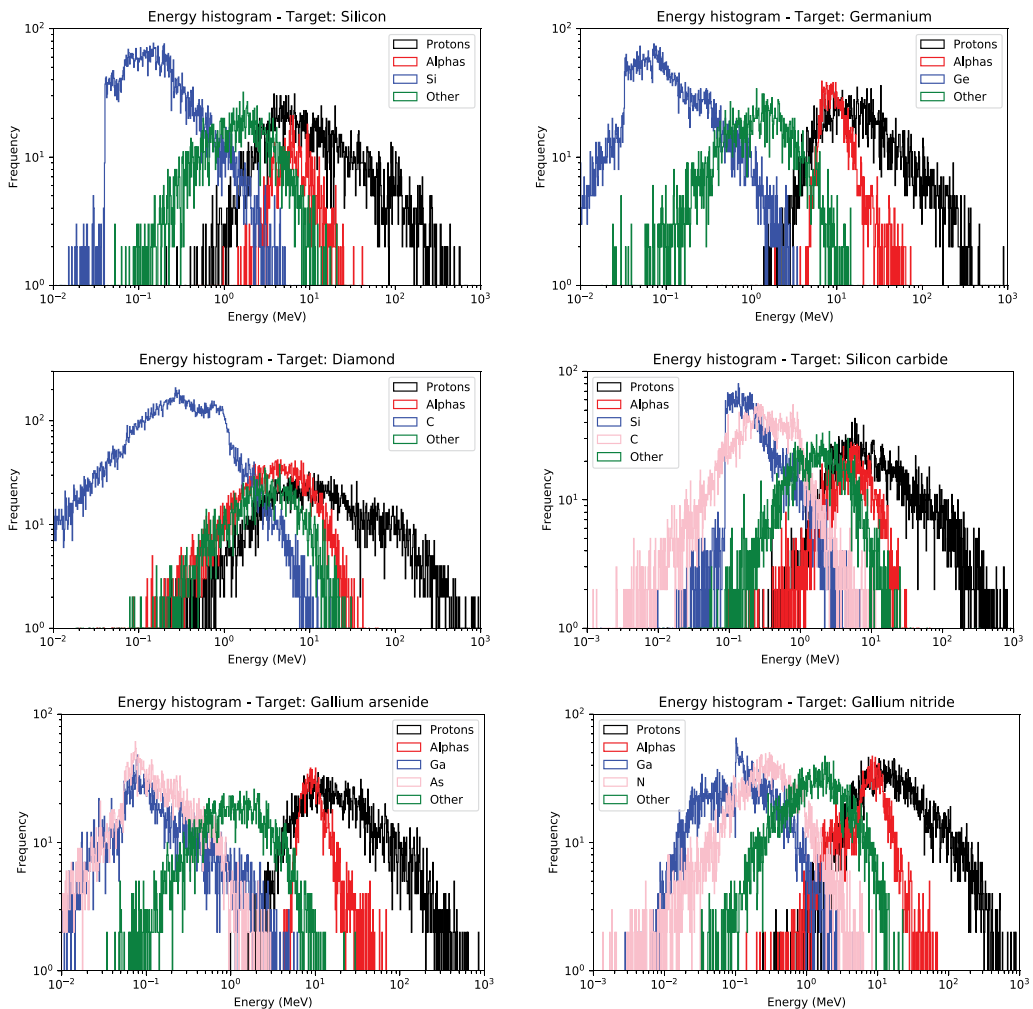


Figure 7. Histograms in energy for recoil nuclei, protons, and alpha particles induced by neutron interactions in the different targets.

LET and range have been obtained from SRIM [12] simulation using numerical functions developed from a behavioral modeling of SRIM Tables [13]. These generalized functions $LET(Z, A, E \text{ target})$ and $Range(Z, A, E \text{ target})$ have been written in C++ and allow us to calculate the two quantities for any given ion defined by the triplet $(Z, A, E \text{ energy})$ and for the semiconductor targets Si, Ge, C, SiC, GaAs, and GaN.

All histograms shown in **Figures 7–10** have been constructed using linear bins with constant bin widths corresponding to the minimum value of the x -axis scale. Curves have been also smoothed in frequency in order to give a more readable aspect to the different distributions.

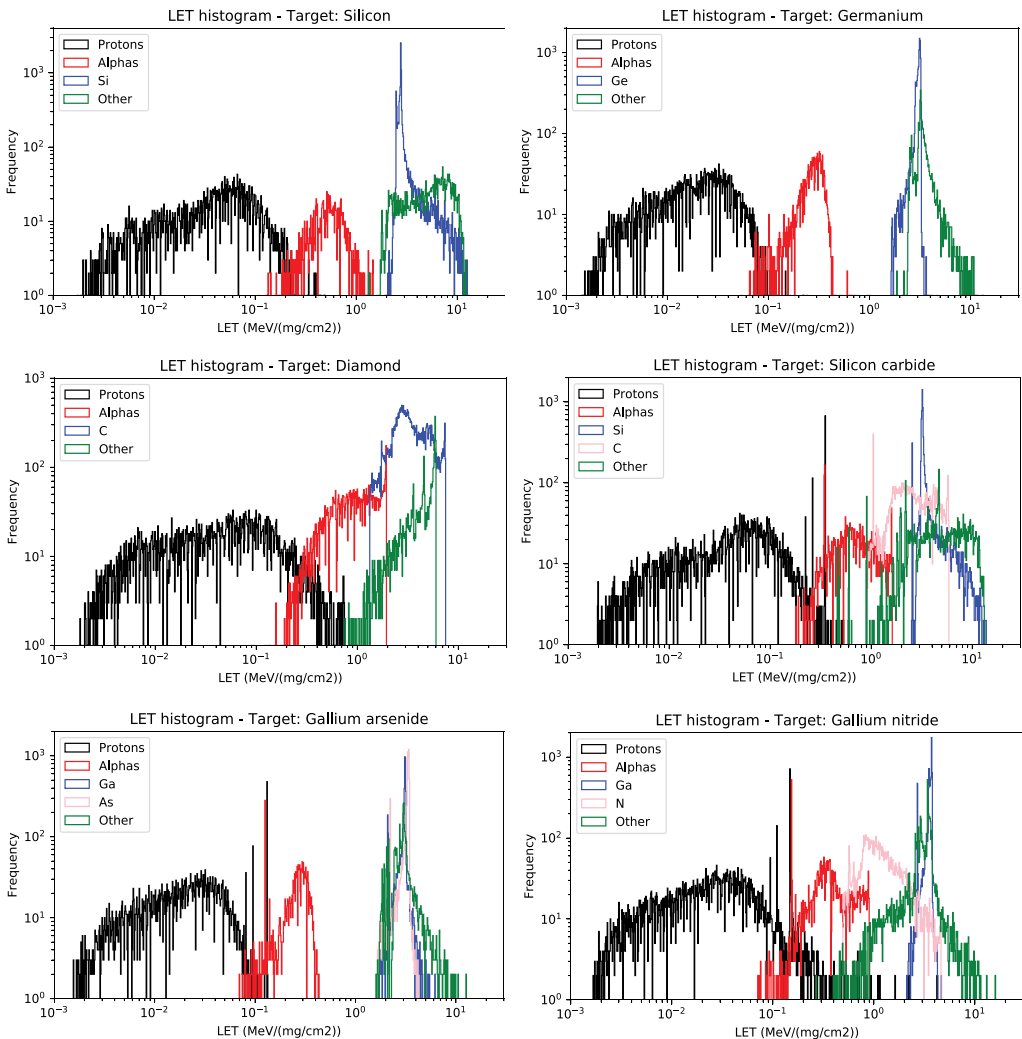


Figure 8. Histograms in LET for recoil nuclei, protons, and alpha particles induced by neutron interactions in the different targets.

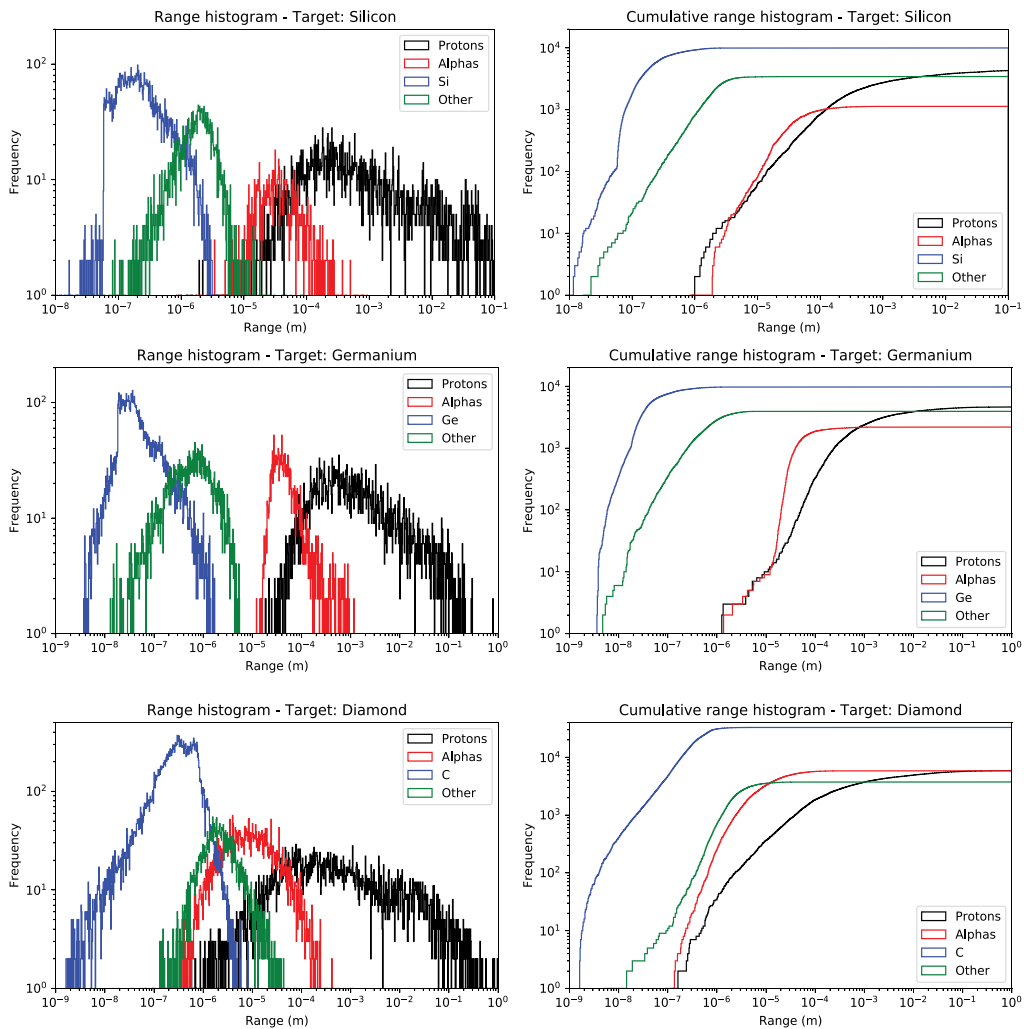


Figure 9. Histograms in range for recoil nuclei, protons, and alpha particles induced by neutron interactions in silicon, germanium, and diamond targets.

Figure 7 shows the histograms in energy for recoil nuclei, protons, and alpha particles induced by neutron interactions in the six targets. We note a certain similarity between the curves related to the different materials:

- Recoil products have systematically a distribution peaking at low energies (1 MeV) and ranging between 1 and 30 MeV for C, 1 and 10 MeV for Si and SiC, 1 to 7–8 MeV for GaAs and GaN, and 1 to 5 MeV for Ge. The high limit in energy of recoil nuclei distributions is much lower than the atomic number Z , which is higher.
- Fragments other than protons and alpha particles show a more broaden distribution than that of recoil products, with a maximum close to 1 MeV and ranging up to a few tens of MeV.

- Alpha particles show a clear peak distribution ranging from 1 MeV to a few tens of MeV and with a maximum around 7–10 MeV, except for C and SiC (around 3 MeV). Above 10 MeV, alpha particles are the most numerous with protons.
- Protons exhibit a large distribution with a maximum around a few MeV (Si, SiC, C) or 10 MeV (Ge, GaAs, GaN) and with a large tail distribution ranging up to several hundreds of MeV (up to 1 GeV for C). Protons clearly dominate in number from 20 MeV to 1 GeV with respect to all other particles.

Figure 8 shows the same data of **Figure 7** but expressed in LET. This transformation has the advantage to show the ionizing power of secondary products just after their release at interaction

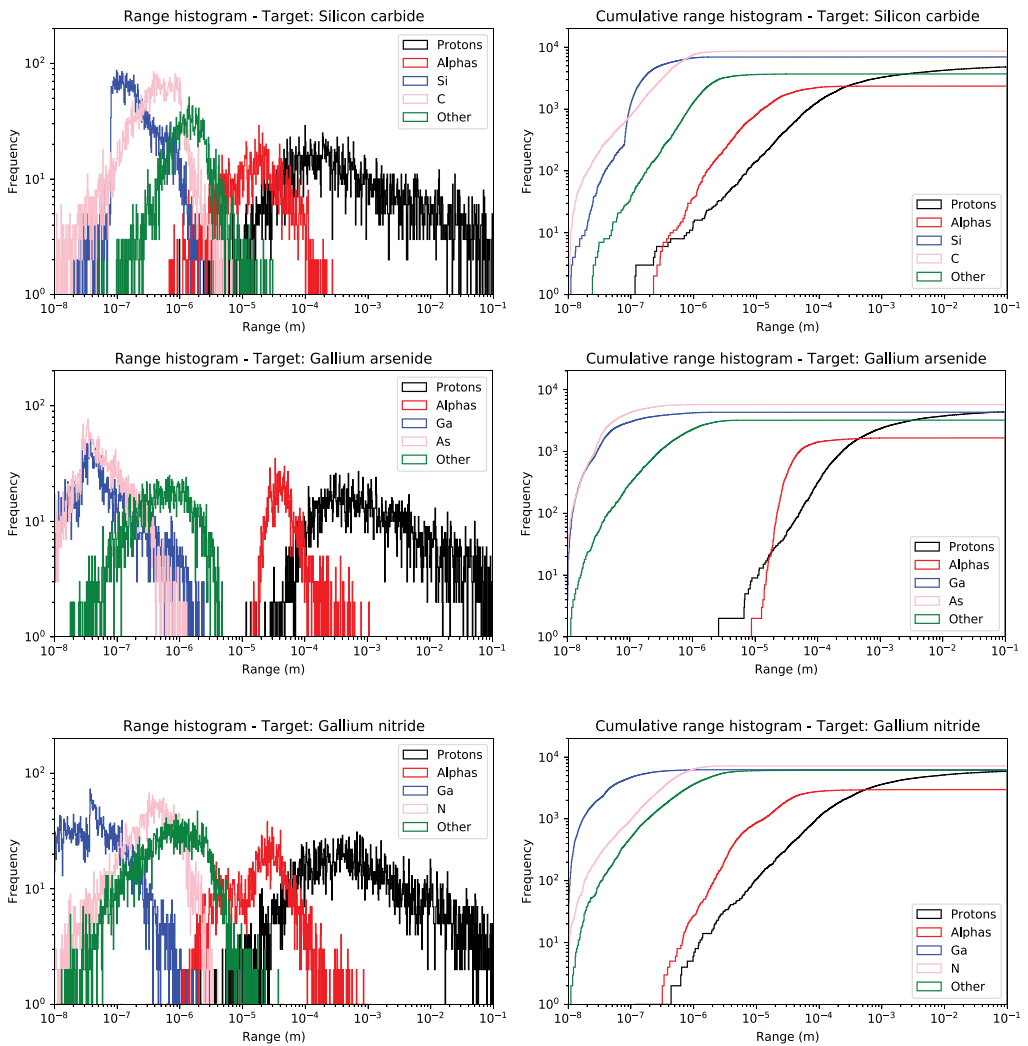


Figure 10. Histograms in range for recoil nuclei, protons, and alpha particles induced by neutron interactions in silicon carbide, gallium arsenide, and gallium nitride targets.

vertex position. In the same way as observed in **Figure 7**, distributions of **Figure 8** show similarities from a target to another:

- Heavy products composed of recoil nuclei and fragments other than protons and alpha particles exhibit the highest LET values, in the range 1 to 10 MeV/(mg/cm²).
- Recoil product LET values are distributed following a very sharp (peaked) distribution centered at approximately 3 MeV/(mg/cm²) for all targets.
- Protons, which are the lightest but the more energetic particles, are characterized by the lowest LET values, from 0.01 to a few 0.1 MeV/(mg/cm²).
- Alpha particles exhibit intermediate LET values, with a peak distribution centered in the range 0.1 to a 1 MeV/(mg/cm²).

To complete previous results, **Figures 9** and **10** show the range distributions of all products, always partitioned in four classes. On the one hand, recoil and other (heavy) products exhibit the shorter ranges, in the deca-nanometer domain and up to a maximum of a few microns. On the other hand, light products, i.e., protons and alpha particles, show much longer ranges up to a few hundreds of microns for the most energetic alpha particles and ranges up to a few millimeters for high-energy protons.

From results shown in **Figures 8–10**, we can logically conclude that recoil nuclei and heavy fragments other than protons and alpha particles are susceptible to induce single events in a very short range from their emission point and with a relative high efficiency, due to their initial high LET values.

On the contrary, protons and alpha particles are characterized by lower LET values but with longer ranges in the different semiconductor materials. Consequently, they are susceptible to induce single events farther from their emission point than heavy fragments up to distances of hundred microns for alpha particles and several millimeters for protons.

3.4. Consequences in terms of electron-hole pair generation and single events

This last paragraph examines the consequences of neutron interactions in the different targets in terms of electron-hole pair generation and fundamental mechanism at the origin of single events in electronics. From the computed databases, we calculated in **Figure 11** the total energy deposited by all the secondary products in the different target materials. Although it is a purely theoretical value, this total amount of deposited energy by ionization process is in the range of 10¹¹ eV for 100 million incident neutrons, which gives an average value in the range of keV per incident neutron, more precisely 1.7 keV for Si, 2.85 keV for Ge, 3.07 keV for C, 2.33 keV for SiC, 2.6 keV for GaAs, and 3.74 keV for GaN. This quantity is found to be minimum for the silicon target and maximum for GaN (**Figure 11**).

Dividing this total energy deposited by all the secondary products by the average energy for creation of an electron-hole pair (given in **Table 1**) gives, for each target material, the upper theoretical limit of the total amount of electron-hole pairs induced by neutrons (via the secondary products). This quantity is shown in **Figure 12** for the different target materials. Also, normalized per incident neutron, this corresponds to 472 e-h pairs for Si, 983 for

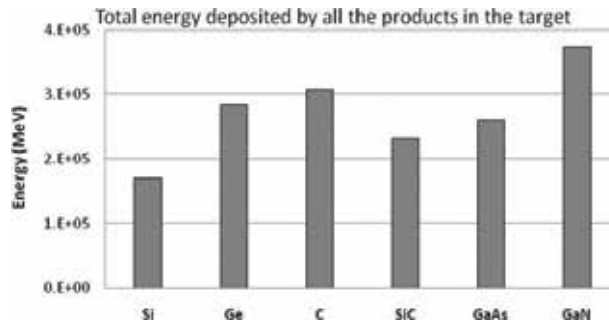


Figure 11. Total energy deposited by all the secondary products induced by neutron interactions in the different target semiconductor materials.

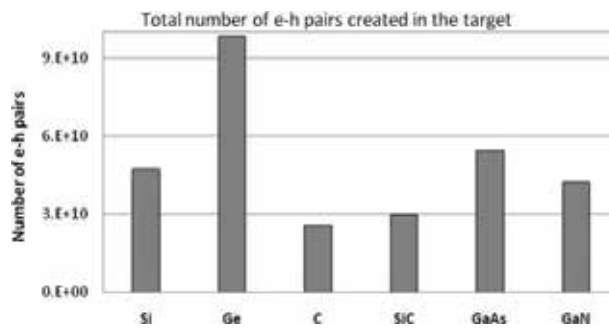


Figure 12. Number of electron-hole pairs created in the different target semiconductor materials by conversion of the total energy deposited by all the secondary products (Figure 11), taking into account the average energy for creation of an electron-hole pair given in Table 1.

Ge, 256 for C, 299 for SiC, 542 for GaAs, and 420 for GaN. These data show that germanium corresponds to the worst case and diamond (also SiC) to the best case with regard to production mechanisms of single events, Si, GaAs, and GaN being relatively equivalent with respect to this criterion. This result can be for the most part explained by the values of the average energy for creation of an electron-hole pair which is very low for Ge (2.9 eV) and extremely important for C (12 eV).

4. Conclusion

In this chapter, we presented a detailed study using extensive Geant4 numerical simulation of nuclear events resulting from the interaction of atmospheric neutrons at the terrestrial level with a target layer composed of various group-IV and III-V semiconductor materials including silicon, germanium, silicon carbide, carbon-diamond, gallium arsenide, and gallium nitride materials. The neutron interaction responses of these different semiconductors have been finely compared in terms of nuclear processes, recoil products, secondary ion production, and fragment energy distributions. Our results show that Si exhibits the

lowest total number of interaction events and, except the particular case of carbon-based materials (diamond and SiC), GaN shows the highest event rate with more than 50% of the supplementary events with respect to Si. Diamond shows a very different behavior than the other materials (our simulation results show a particularly elevated number of elastic events for diamond as compared to other materials) since it is an excellent neutron moderator. For silicon carbide, which can be viewed as a “mixture” of Si and C at the atomic level, it shows an intermediate behavior between Si and C with a number of elastic events quasi $\times 2$ with respect to Si due to the presence of C. Concerning the fraction of elastic and inelastic events, three different behaviors can be highlighted: low ($< 30\%$, Ge and GaAs), intermediate ($40\text{--}60\%$, Si and GaN), and high ($> 60\%$, C and SiC) elastic event rates; the presence of low-Z elements such as C and N, respectively, in SiC and GaN leads to increase elastic interactions in these last materials with respect to the elastic rates observed for Si and GaAs. For Si, Ge, and GaAs, the total number of generated secondary products is in the same order of magnitude (with $\text{Si} < \text{Ge} < \text{GaAs}$). For GaN and SiC, the presence of low-Z elements increases the number of elastic events and indirectly increases the total number of secondary products. Carbon-diamond shows the highest number of events/products due to its high power of neutron moderation. Concerning the nature of the secondary products, our simulations show that, for the six materials, the most frequent produced secondary particles are the recoil products due to neutron elastic interactions with the nuclei of the semiconductor lattice, followed by protons in the second position and alpha particles in the third one. All the other products are systematically less produced than these three categories of products. A detailed analysis of the secondary ions produced during neutron interactions in the different target materials in terms of initial energy (when products are released), linear energy transfer (LET), and range in the target material has been also conducted. Recoil nuclei and heavy fragments have been shown to be susceptible to induce single events in a very short range from their emission point and with a relative high efficiency, due to their initial high LET values. On the contrary, protons and alpha particles, characterized by lower LET values but longer ranges in the different semiconductor materials, are susceptible to induce single events farther from their emission point than heavy fragments up to distances of hundred microns for alpha particles and several millimeters for protons. Finally, the consequences of neutron interactions in the different targets in terms of electron–hole pair generation, a fundamental mechanism at the origin of single events in electronics, have been examined. Our results show that germanium corresponds to the worst case and diamond (also SiC) to the best case with regard to e–h pair production, Si, GaAs, and GaN being relatively equivalent and of intermediate behavior with respect to this criterion.

Author details

Daniela Munteanu and Jean-Luc Autran*

*Address all correspondence to: jean-luc.autran@univ-amu.fr

Aix Marseille Université, CNRS, Université de Toulon, Marseille, France

References

- [1] International Technology Roadmap for Semiconductors 2.0. Online available: <http://www.itrs2.net>
- [2] Serre S, Semikh S, Uznanski S, Autran J-L, Munteanu D, et al. Geant4 analysis of n-Si nuclear reactions from different sources of neutrons and its implication on soft-error rate. *IEEE Transactions on Nuclear Science*. 2012;**59**(4, 1):714-722
- [3] Liu H, Cotter M, Datta S, Narayanan V. Technology Assessment of Si and III-V FinFETs and III-V Tunnel FETs from Soft Error Rate Perspective. International Electron Device Meeting, Institute of Electrical and Electronics Engineers. 2012. p. 577-580
- [4] Agostinelli S, et al. Geant4 – a simulation toolkit, *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment*. 2003;**506**(3):250-303
- [5] Allison J, et al. Recent developments in Geant4, *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment*. 2016;**835**:186-225
- [6] Madelung O. *Semiconductors Data Handbook*. Berlin Heidelberg; Springer-Verlag. 2004. 691 p
- [7] Gordon MS, Goldhagen P, Rodbell KP, Zabel TH, Tang HHK, Clem JM, Bailey P. Measurement of the Flux and Energy Spectrum of Cosmic-Ray Induced Neutrons on the Ground. *IEEE Transactions on Nuclear Science*. 2004;**51**:3427-3434
- [8] Autran JL, Munteanu D. *Soft Errors: from particles to circuits*. Taylor & Francis/CRC Press. 2015. 439 p
- [9] Geant4 General Particle Source (GPS). Online available: <https://geant4.web.cern.ch/geant4/UserDocumentation/UsersGuides/ForApplicationDeveloper/html/ch02s07.html>
- [10] Geant4 version 4.9.4. Online available: http://geant4.cern.ch/collaboration/working_groups/electromagnetic/physlist9.4.shtml
- [11] Wrobel F, Gasiot J, Saigné F, Touboul AD. Effects of atmospheric neutrons and natural contamination on advanced microelectronic memories. *Applied Physics Letters*. 2008. p. 064105
- [12] Ziegler JF, Ziegler MD, Biersack JP. SRIM – The stopping and range of ions in matter (2010). *Nuclear Instruments and Methods in Physics Research Section B: Beam Interactions with Materials and Atoms*. 2010;**268**(11):1818-1823
- [13] Martinie S, Saad-Saoud T, Moindjie S, Munteanu D, Autran JL. Behavioral modeling of SRIM tables for numerical simulation, *Nuclear Instruments and Methods in Physics Research Section B: Beam Interactions with Materials and Atoms*. 2014;**322**:2-6

Ultrashort Pulse Generation in Ce:LiCAF Ultraviolet Laser

Marilou Cadatal-Raduban, Minh Hong Pham,
Luong Viet Mui, Nguyen Dai Hung and
Nobuhiko Sarukura

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.73501>

Abstract

Transient cavity method used to generate ultrashort laser pulses in dye lasers is extended to a solid-state gain medium. Numerical simulations are performed to investigate the spectro-temporal evolution of broadband ultraviolet (UV) laser emission from Ce³⁺-doped LiCaAlF₆ (Ce:LiCAF), which is represented as a system of two homogeneous broadened singlet states. By solving the rate equations extended to multiple wavelengths, the appropriate cavity length and Q-factor for optimal photon cavity decay time and pumping energy that will generate resonator transients is determined. Formation of resonator transients could generate picosecond UV laser pulses from a Ce:LiCAF crystal pumped by the fourth harmonics (266 nm) of a Nd:YAG laser. Numerical simulations indicate that a 1-mol% Ce³⁺-doped LiCAF crystal that is 1-mm long can generate a single picosecond pulse. This is accomplished by using a low Q (output coupler reflectivity of 10%), short cavity (cavity length of 2 mm) laser oscillator. Ultrashort pulses can also be generated using other rare earth-doped fluoride laser materials using this technique.

Keywords: numerical simulation, ultrashort pulse, ultraviolet, resonator transient, Ce:LiCAF crystal

1. Introduction

Tunable ultrashort-pulsed laser emission in the ultraviolet (UV) region is highly sought after because of its numerous applications in many fields of science and technology [1, 2]. Ultrashort pulses are necessary for controlling ultrafast chemical processes [1], probing fast physical and chemical processes [3, 4], and investigating the relaxation of charge carriers in conductors [5],

to name a few. On the other hand, UV lasers have many applications in various fields such as surface structuring [6–8], micromachining [9], remote sensing [10], spectroscopy and imaging [11]. An important advantage of having ultrashort pulses in the UV wavelength region is its ability to modify material properties only within the laser focus where the peak power is high. This feature is especially critical in micromachining [12]. Ultrashort UV pulses also permit outstanding temporal resolutions for pump-probe experiments [13]. Available light sources that satisfy these requirements are limited, despite the many applications. Excimer lasers can emit UV wavelengths, but these are bulky and cumbersome to maintain [14, 15]. UV laser emission using frequency conversion in nonlinear crystals is well established but is complex, has limited spectral bandwidth, non-tunable, and has low conversion efficiency. In contrast, high-power, all-solid-state UV lasers are highly regarded for simplicity in operation and maintenance. Hence, there is a great deal of interest for developing all-solid-state UV lasers.

Cerium ion (Ce^{3+})-doped wide band gap fluorides have been the most successful tunable solid-state laser media in the UV region. Direct UV emission has been reported from Ce^{3+} -doped YLiF_4 , LaF_3 , LiLuF_4 , LiCaAlF_6 , and LiSrAlF_6 crystals [16–23]. Among these known UV laser crystals, Ce^{3+} -doped lithium calcium hexafluoroaluminate ($\text{Ce}^{3+}:\text{LiCaAlF}_6$ or $\text{Ce}:\text{LiCAF}$) is the most prominent and successful solid-state gain medium for amplifying short UV pulses as well as for generating ultrashort UV pulses because it is highly transparent, has low excited state absorption, not prone to color center formation and therefore tolerant to laser-induced damage. Most importantly, it is absorbing at around 266 nm, which makes the fourth harmonics of readily-available Nd:YAG lasers an ideal excitation source. It also has sufficiently large effective gain cross-section of $6.0 \times 10^{-18} \text{ cm}^2$ that is favorable for oscillators, and a high saturation fluence of 115 mJ/cm^2 . Lastly, it has broad tunability from 280 to 325 nm and enough bandwidth to generate 3-fs pulses [2, 20–22, 24–30].

Direct generation of tunable UV short-pulses using solid-state gain media, such as $\text{Ce}:\text{LiCAF}$, is not as straight forward as using near IR tunable solid-state laser media, such as $\text{Ti}:\text{sapphire}$, due to the difficulty of obtaining continuous wave (CW) laser operation. As a result, Kerr lens mode-locking schemes that utilize spatial or temporal Kerr type nonlinearity are a real challenge in the UV region [31]. The lifetime of the upper energy level of $\text{Ce}:\text{LiCAF}$ is about 25 ns, which is too short to directly generate tunable UV short-pulses. Therefore, high pump power densities are required to achieve CW and mode-locked operation [32]. Moreover, the cavity lengths of the pump and the $\text{Ce}:\text{LiCAF}$ laser oscillator have to be matched for synchronous mode locking [33]. A typical oscillator in mode-locking schemes also uses a four-mirror z-fold cavity, which means that reducing losses in the laser cavity is also crucial. On the other hand, the transient cavity method is a simpler means of generating tunable UV ultrashort pulses because it only utilizes the usual two-mirror laser oscillator. Pulse shortening by resonator transients has been demonstrated in dye lasers where laser pulse durations that are an order of magnitude shorter than the pump pulse have been obtained [34–36]. This book chapter discusses numerical simulations that extend the resonator transient technique to solid-state gain media. Since formation of resonator transients that lead to picosecond pulses depend on the cavity length and the Q-factor of the laser oscillator cavity, numerical simulations are carried out to investigate the decay time of the photons within the laser cavity as well as the energy of the picosecond Nd:YAG pump laser that will support the formation of these resonator transients. The technique can be extended to other rare earth-doped fluoride laser materials.

2. Review of cerium ion-doped fluoride crystals

Laser gain media based on wide bandgap fluoride hosts was first proposed as a result of spectroscopic studies on trivalent lanthanides such as neodymium (Nd), cerium (Ce), and thulium (Tm) doped in solid-state hosts [37]. The intense broad band UV fluorescence from 276 to 312 nm observed from $\text{Ce}^{3+}:\text{LaF}_3$ and the 288–322 nm fluorescence from $\text{Ce}^{3+}:\text{LuF}_3$ were attributed to dipole-allowed 5d to 4f (5d-4f) radiative transition of Ce ions in the LaF_3 and LuF_3 fluoride hosts [38]. Subsequently, the first laser emission from the 5d-4f transition was achieved in 1977 using $\text{Ce}^{3+}:\text{YLiF}_4$. It emitted at 325.5 nm when it was optically pumped at 249 nm [16]. However, the progress of $\text{Ce}^{3+}:\text{YLiF}_4$ was limited by poor performance characteristics brought about by an early onset of saturation and roll off in the above-threshold gain and power output as well as a drop in the output for pulse repetition rates above 0.5 Hz. Although lasing from $\text{Ce}^{3+}:\text{YLiF}_4$ is ground breaking, the existence of solarization or color center formation prevents this material from being of practical use, thereby hindering its further development. In 1980, operation of an optically pumped $\text{Ce}^{3+}:\text{LaF}_3$ laser was reported [17]. Limitations of this laser medium include low output power and high lasing threshold. Moreover, the lasing results have not been reproduced. Subsequent experiments on other Ce^{3+} -doped fluorides have not been very successful due to the formation of transitory or permanent color centers. Such color centers were essentially due to absorption of the pump and/or the laser radiation from emitting 5d states leading to the promotion of an electron into the conduction band followed by trapping by impurities or defects [39–43]. Recent investigations, however, showed that by an appropriate choice of activator-matrix complexes and active medium-pump source combinations, efficient tunable lasers using d-f transitions can be created [19, 21, 22, 44]. In 1992, emission from $\text{Ce}^{3+}:\text{LiLuF}_4$ pumped by a KrF excimer laser was reported. This laser material has almost the same optical properties as $\text{Ce}^{3+}:\text{YLiF}_4$ but with smaller solarization effect [19]. As a result, slope efficiencies of more than 50% were obtained [45, 46]. Moreover, continuous tunability was achieved from 305 to 333 nm [46]. Subsequently, lasing from $\text{Ce}^{3+}:\text{LiCaAlF}_6$ (Ce:LiCAF) was reported. This was a milestone not only because Ce:LiCAF can be pumped by the fourth harmonic of a Nd:YAG laser, but also because remarkably, no solarization effect was observed

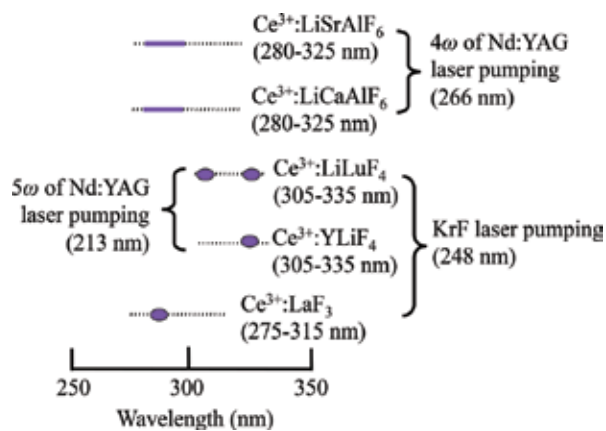


Figure 1. Ce^{3+} -doped lasers for tunable UV radiation. Solid lines and dots indicate the confirmed tunable wavelength region, dotted lines show potential tunable wavelength region.

in this crystal [21, 22]. Following the success of Ce:LiCAF, lasing from Ce³⁺:LiSrAlF₆ was reported. It can also be pumped with the fourth harmonic of a Nd:YAG laser and it has similar laser properties as the Ce:LiCAF crystal [23, 47]. **Figure 1** summarizes the tunable wavelength regions of the five currently known Ce-doped lasers.

Ce:LiCAF has the following advantages over the other known Ce³⁺-doped fluoride hosts: (1) the strong absorption band at 266 nm would allow direct optical pumping by the fourth harmonics of the Nd:YAG laser; (2) the wide fluorescence band offers a tuning range from the 280 to 320 nm for possible pulse compression; (3) the gain cross-section of Ce:LiCAF ($6 \times 10^{-18} \text{ cm}^2$) is high, even higher compared to the Ti:Sapphire. This property is ideal for building a laser resonator; (4) the saturation fluence ($\sim 115 \text{ mJ/cm}^2$) of Ce:LiCAF is higher than organic dyes. This is significant when using this crystal in power amplifiers; (5) the nanosecond lifetime of Ce:LiCAF maybe too short to for regenerative amplifications, but this is sufficiently long to allow multi-pass amplification.

3. Cerium ion dopant

Figure 2 shows the energy level structure of Ce³⁺ doped into a fluoride host. The ground state 4f configuration has two energy levels, ²F_{5/2} and ²F_{7/2}, which are separated by 0.2793 eV (2253 cm⁻¹). The excited state 5d configuration also has two energy levels, ²D_{3/2} and ²D_{5/2}, which are separated by 6.166 eV (49,733 cm⁻¹) and 6.475 eV (52,226 cm⁻¹) from the ground state, respectively. The exact positions of these energy levels would depend on the specific host. Lasing in the UV is based on the electric dipole-allowed interconfigurational 5d-4f transitions. In contrast, conventional trivalent lanthanide laser crystals, such as Nd:YAG, uses the intraconfigurational 4f-4f transition that results to infrared (IR) emission. As a result, UV

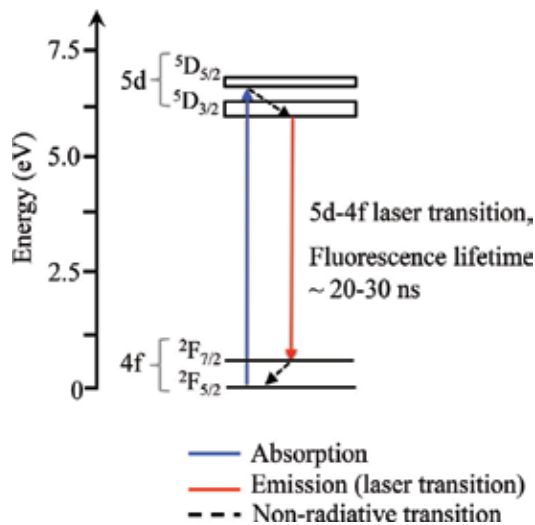


Figure 2. Energy level structure of Ce³⁺ doped into a fluoride host.

fluorescence emission from Ce^{3+} -doped fluoride crystals have smaller radiative lifetimes of a few tens of nanoseconds compared to IR emissions that have lifetimes within hundreds of microseconds. In addition, fluorescence from the 5d-4f transitions is characterized by broad bandwidths and large Stokes shifts. The broad gain bandwidth enables tunability and ultra-short laser pulse generation from Ce^{3+} -activated laser crystals. The large energy gap between the excited state 5d configuration and the 4f ground state configuration laser levels results to low multi-phonon related nonradiative decay. Therefore, quantum efficiencies as high as 90% are expected from Ce^{3+} -activated laser crystals [20, 48].

4. Numerical simulation of the electronic properties of the LiCAF host

LiCAF is a colquiriite-type fluoride with a hexagonal crystal structure belonging to the P-31c space group (group number 163). It is optically a uniaxial crystal with two formula units per unit cell. Six fluorine (F) atoms surround a lithium (Li), calcium (Ca), or aluminum (Al) atom. Each Li, Ca, and Al cation occupies a deformed octahedral site as shown in **Figure 3a**. This structure is also described by an alternative stacking of metallic and fluorine atom layers parallel to the c-axis [49–51]. The fraction coordinates of the representative atoms in the unit cell are shown in **Table 1**.

First-principles density functional theory (DFT) calculations are used to obtain the optimized volume, electronic band structure, total and partial density of states (DOS), and the band gap

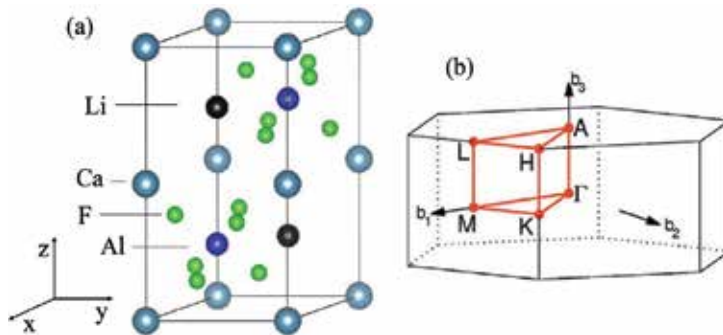


Figure 3. (a) Colquiriite-type structure of LiCAF and (b) first Brillouin zone of the hexagonal unit cell of a LiCAF crystal.

	x	y	z
Li	1/3	2/3	1/4
Ca	0	0	0
Al	2/3	1/3	1/4
F	0.3769	0.0312	0.1435

Table 1. Atomic positions of the LiCAF atoms.

energies of the LiCAF crystal. These calculations employed the projector-augmented wave (PAW) method as implemented within the Vienna Ab Initio Simulation Package (VASP) [52–57], with a plane-wave basis cutoff of 500 eV and a hybrid density functional, which uses the full Perdew-Burke-Ernzerhof (PBE) [58, 59] correlation energy but mixes 65% PBE exchange with 35% exact exchange [60–63]. The initial charge density and wave function was generated using a $3 \times 3 \times 1$ Monkhorst-Pack k-point grid. For the band structure and DOS diagrams, the k-points were chosen following the first Brillouin zone and the path $\Gamma \rightarrow M \rightarrow K \rightarrow A \rightarrow L \rightarrow H$ shown in **Figure 3b** [64].

The electronic band structure of LiCAF along the high symmetry lines of the first Brillouin zone is shown in **Figure 4**. The maximum of the valence band is located at the k-point between M and K while the conduction band minima is at the Γ point. Therefore, LiCAF has an indirect band gap with a band gap energy of 12.23 eV. This result is 3.30% and 10.51% different from the experimentally obtained band gap energies of 12.65 eV and 11.07 eV, respectively [65–67]. **Figure 5** shows the total and partial DOS of LiCAF. The maximum valence band is derived from the fluorine 2p states whereas the aluminum 4s and fluorine 3s states contribute to the minimum conduction band.

Excited state absorption (ESA), which is prevalent in rare-earth-doped fluorides operating in the UV region, has been observed in both $\text{Ce}^{3+}:\text{LiCaAlF}_6$ (Ce:LiCAF) and $\text{Ce}^{3+}:\text{LiSrAlF}_6$ (Ce:LiSAF). However, experimental investigations reveal that Ce:LiSAF experiences ESA to a greater extent compared to Ce:LiCAF and therefore, the conversion efficiency of a Ce:LiSAF laser is lower. ESA results from an electron being promoted from the 5d excited state configuration of Ce^{3+} to the conduction band of the LiCAF host [24, 68]. Therefore, the onset of ESA strongly depends on the host. If the conduction band minimum of the host is close to the 5d excited state level of the activator ion, then ESA will be greater in this laser material. Similar band structure and DOS calculations performed for LiSAF reveal that the band gap energy of LiSAF is 11.79 eV, which is 0.44 eV smaller than LiCAF [62]. Associated with the strong ESA is color center formation or solarization, which happens due to an electron getting trapped at impurities in the conduction band of the host as shown in **Figure 6** [27]. Under UV excitation, color centers can be created due to solarization. Broad absorption bands in energies other than the band gap then appear as a result of color center-formation. The Ce^{3+} ions that are doped in

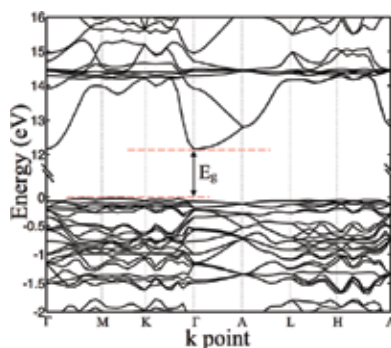


Figure 4. Simulated electronic band structure of LiCAF host crystal.

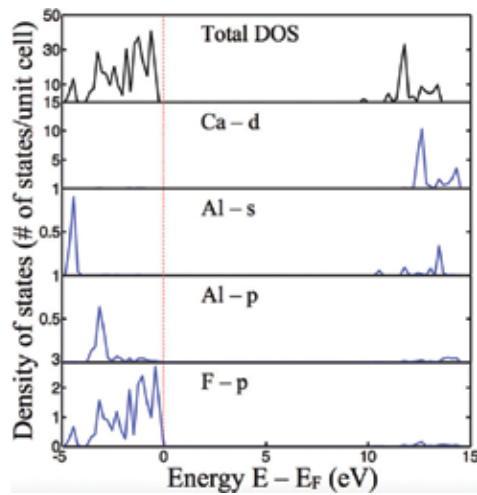


Figure 5. Total and projected density of states of LiCAF host.

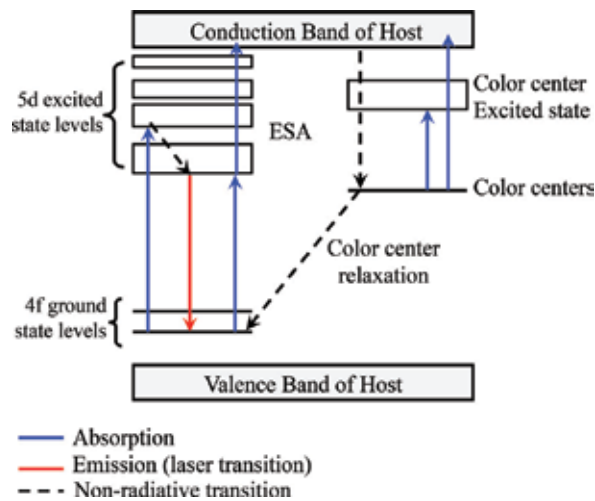


Figure 6. Excited state absorption (ESA) and color center formation in Ce³⁺-doped fluoride gain media [27].

the LiCAF or LiSAF host tend to occupy the Ca²⁺ or Sr²⁺ octahedral sites [29]. The CeF₆ cluster in LiSAF will then have to compensate for the disparity in the size of the Ce³⁺ ion, which is 1.15 Å, and the Sr²⁺ ion, which is 1.27 Å. This compensation makes LiSAF prone to defects such as cracks and impurities during crystal growth [69]. For comparison, the size of the Ca²⁺ ion in LiCAF, which is 1.14 Å, is similar to that of the Ce³⁺ ion. Ce:LiCAF does not exhibit as much defects as Ce:LiSAF under the same growth conditions [69]. The larger amount of cracks and defect in Ce:LiSAF as well as the significant ESA as described above result to more pronounced solarization in this crystal. Both ESA and solarization compete with the lasing process and these results to lower laser conversion efficiency.

5. Numerical simulation of optimized resonator transients for generating ultrashort Ce:LiCAF laser pulses

Sub-nanosecond (0.7 ns), single-pulse, ~290 nm laser emission from a Ce:LiCAF crystal using an oscillator with a 25 mm cavity length and 25% output coupler reflection has been demonstrated experimentally. In this report, the end-pumping configuration was used to excite the crystal with an Nd:YAG laser emitting 266-nm wavelength (fourth harmonics), 5-ns pulses at 10 Hz repetition rate [70]. Another study reported the generation of 150-ps laser pulses from a 10-mm long Ce:LiCAF crystal. In this report, the laser oscillator was established using a cavity length of 15 mm, output coupler reflection of 30%, and 75 ps excitation pulses [71]. Despite the short pulse durations reported in these studies, the inherent properties of Ce:LiCAF provides the capability for this solid-state laser gain medium to generate even shorter UV pulse durations, which up to this point has not been fully realized. Numerical calculations play an important role in determining the influence of optical parameters, such as pumping energy, Q-value and cavity length on the output energy and pulse duration of a Ce:LiCAF UV laser, and hence in optimizing the laser oscillator design.

The Ce:LiCAF crystal used in the numerical calculations has 1 mol% Ce³⁺ ion doping and is 1 cm long. This crystal is placed inside a Fabry-Perot laser cavity of length L . It is end-pumped by the fourth harmonics (266 nm) of a ps Nd:YAG laser. The pump pulse is a Gaussian beam with 75 ps (FWHM) pulse duration. The end mirror of the laser cavity is flat with reflectivity denoted by R_1 . The output coupler of the laser cavity is also flat with reflectivity denoted by R_2 . It is assumed that both mirrors have uniform reflectivity within the emission bandwidth of the laser crystal. The optical properties of the Ce:LiCAF crystal given in **Table 2**, which is detailed in several papers, are used as calculation parameters [2, 20–22, 24–30].

Laser emission is approximated using a system of two homogeneous broadened singlet states. Eqs. (1)–(4) show the modified rate equations as it applies to multiple wavelengths, which accurately simulates the broad emission bandwidth of the Ce:LiCAF UV laser [72–74]:

$$\frac{\partial N_1}{\partial t} = P(t) + \left[\sum_1^n \sigma_{ai} I_i \right] N_0 - \left[\sum_1^n \sigma_{ei} I_i + \frac{1}{\tau} \right] N_1 \quad (1)$$

$$P(t) = \frac{P_{in} [1 - \exp(-\alpha l)] \lambda_p}{hc\tau r^2 l} \exp \left[\frac{-4 \ln(2)(t - t_0)^2}{\Delta t^2} \right] \quad (2)$$

$$\frac{\partial I_i}{\partial t} = [2(\sigma_{ei} N_1 - \sigma_{ai} N_0) l - \beta] \frac{I_i}{T} + A_i N_1 \quad (3)$$

$$T = \frac{2[L + l(n - 1)]}{c} \quad (4)$$

The population density in the upper laser state as a function of time is given by Eq. (1) where N_0 is the lower-state population density, N_1 is the upper-state population density, $N = N_0 + N_1$ is the total doping density, I_i is the intensity of the laser with wavelength λ_i , τ is the fluorescence decay

Ce ³⁺ doping concentration	1 mol%
Doping density, N	$5 \times 10^{17} \text{ cm}^{-3}$
Absorption coefficient at 266 nm (wavelength of pump laser)	4 cm^{-1}
Absorption cross-section, σ_{ai}	$2.606 \times 10^{-19} \text{ cm}^{-2}$ at 290 nm
Emission cross-section, σ_{ei}	$9.6 \times 10^{-18} \text{ cm}^{-2}$ at 290 nm
Refractive index, n	1.41
Fluorescence lifetime, τ	25 ns
Spontaneous emission constant, A_i	$0.2 \times 10^{-10} \text{ cm s}^{-2}$
Wavelength of pump laser, λ_p	266 nm
Radius of pump beam inside the crystal	100 μm
Pulse duration of pump pulse	75 ps
Planck's constant, h	$6.62606957 \times 10^{-34} \text{ J s}$
Speed of light, c	$3 \times 10^8 \text{ m s}^{-2}$
Reflectivity of end mirror, R_1	100%

Table 2. Optical properties of the Ce:LiCAF crystal and values of parameters that were kept constant in numerical simulations.

time. P is the rate of pumping which is further described by Eq. (2) where λ_p is the pump laser's wavelength, P_{in} is the power of the pump, l is the Ce:LiCAF crystal's length, r is the pump beam's radius inside the laser medium, h is Planck's constant which is $6.62606957 \times 10^{-34} \text{ J s}$, c is the speed of light which is $3 \times 10^8 \text{ m s}^{-1}$, α is the absorption coefficient of Ce:LiCAF at 266 nm (pump laser's wavelength), t is the duration of pumping, Δt is the laser pump's pulse duration, and t_0 is the time when P_{in} is maximum. P has a unit of s^{-1} . The laser intensity inside the cavity as a function of time is given by Eq. (3) where β is the round-trip loss defined by $\beta = -\ln(R_1R_2)$ in which case R_1 is the end mirror's reflectivity and R_2 is the output coupler's reflectivity, σ_{ei} is the emission cross-section at wavelength λ_i , σ_{ai} is the absorption cross-section at wavelength λ_i [24], A_i is a constant that simulates spontaneous emission at wavelength λ_i and its value is considered equal for all wavelengths since the duration of pumping is much greater than the memory time of the system of equations [72–74]. In this work, the value of A_i is about three times longer than the decay time of fluorescent dyes as estimated using the fluorescence decay time of Ce:LiCAF, which is about 25 ns. T is the cavity round-trip time as defined by Eq. (4) where L is the laser cavity's length and n is the Ce:LiCAF crystal's refractive index. The values of the constants used in the simulations are given in **Table 2**.

Rate Eqs. (1)–(4) were used to model the experimental results [71, 75] by using the same parameters that were used in the experiment [71]. The Ce:LiCAF crystal length, $l = 10 \text{ mm}$; $N = 5 \times 10^{17} \text{ cm}^{-3}$ for a doping concentration of 1 mol%; refractive index, $n = 1.41$; length of cavity, $L = 15 \text{ mm}$; reflectivity of end mirror, $R_1 = 100\%$; reflectivity of output coupler, $R_2 = 30\%$; radius of pump beam, $r = 100 \mu\text{m}$; energy of pump, $E_{\text{pump}} = 94 \mu\text{J}$; and the laser pump's pulse duration, $\Delta t = 75 \text{ ps}$. These values are given in **Tables 2** and **3**. The spectro-temporal plot of the simulated broadband emission from 286 to 290 nm is shown in **Figure 7**. The corresponding

spectral dynamics derived from integrating along the horizontal axis of **Figure 7** is shown in **Figure 8**. The maximum laser intensity is observed at around 288.5 nm. This corresponds to the wavelength where the gain coefficient is also maximum. The distinct feature observed at around 289.5 nm is consistent with experimental observation [20–30] and is reported to be present for

Crystal length, l	10 mm
Cavity length, L	15 mm
Reflectivity of output coupler, R_2	30%
Energy of pump beam	94 μ J

Table 3. Ce:LiCAF resonator parameters used to reproduce experimental results.

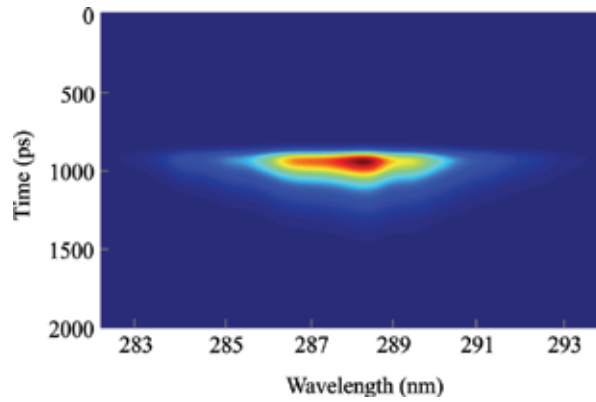


Figure 7. Spectro-temporal evolution of the low-Q, short cavity Ce:LiCAF laser pulse emission obtained numerically using the rate equations. The oscillator parameters are $N = 5 \times 10^{17} \text{ cm}^{-3}$, $l = 10 \text{ mm}$, $L = 15 \text{ mm}$, $R_1 = 100\%$, $R_2 = 30\%$, $r = 100 \mu\text{m}$, $\Delta t = 75 \text{ ps}$, $E_{\text{pump}} = 94 \mu\text{J}$, and $n = 1.41$.

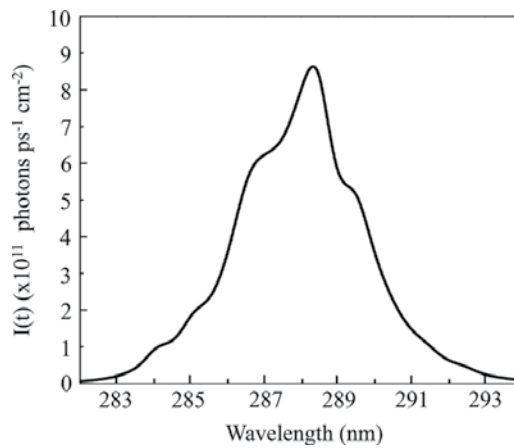


Figure 8. Spectral dynamics of the numerically calculated spectro-temporal profile in **Figure 7**.

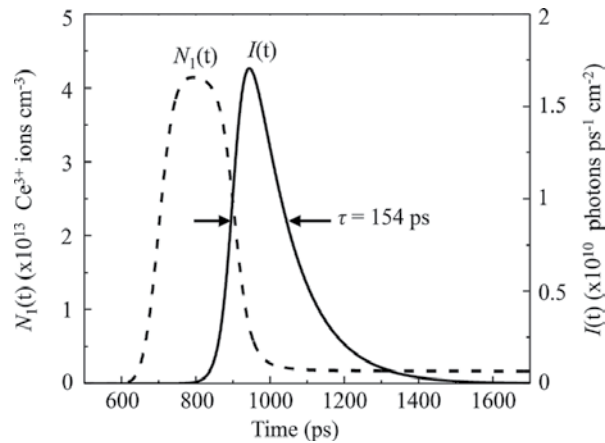


Figure 9. Temporal dynamics of the numerically calculated spectro-temporal profile in **Figure 7**.

any amount of doping. **Figure 9** shows the temporal dynamics of the laser emission, which was derived by integrating along the vertical axis of **Figure 7**. Electron density builds up in the upper laser state of Ce:LiCAF until population inversion is achieved. Lasing threshold is reached around 0.3 ns after the onset of the pump pulse (not shown). The lasing threshold was measured when laser emission is observed. Peak emission is achieved about 0.5 ns after the onset of the pump pulse. The pulse duration is estimated from the full-width-at-half-maximum (FWHM) to be around 154 ps. Results of the temporal dynamics simulation are comparable, within the limits of jitter, to experimental results estimated from the streak camera image that was obtained experimentally [71, 75]. The experimental pulse duration is about 150 ps [71]. The good agreement between the numerical and experimental results indicate that the system of two homogeneous broadened singlet states and the modified rate equations as it applies to multiple wavelengths provide good approximations for predicting the experimental outcome. In succeeding calculations, the model described above was then used to optimize the optical parameters in order to achieve ultrashort pulse emission from a Ce:LiCAF solid-state UV laser.

Eqs. (1) and (2) clearly show that the rate of change of excited state population density and laser intensity strongly depend on the pump energy. Therefore, we solved the rate equations for varying pump energies from 25 to 600 μJ in order to determine its effect on the temporal evolution of the Ce:LiCAF laser emission. Results are shown in **Figure 10a–f**. The parameters used in the numerical simulation were kept constant except for the pump energy. These parameters were the same as the experimental parameters, i.e. crystal length, $l = 10 \text{ mm}$; $N = 5 \times 10^{17} \text{ cm}^{-3}$ for 1 mol% doping concentration; cavity length, $L = 15 \text{ mm}$; end mirror reflectivity, $R_1 = 100\%$; output coupler reflectivity, $R_2 = 30\%$; pump beam radius, $r = 100 \mu\text{m}$; and pulse duration of the laser pump, $\Delta t = 75 \text{ ps}$. Lasing is achieved when the pump energy is around 37.5 μJ (**Figure 10b**). **Figure 10c–f** clearly shows that the laser pulse duration shortens as the pump energy is increased up to 360 μJ (**Figure 10d**). As the pump energy is increased, the time delay between the onset of laser emission and the pump also decreases due to the earlier occurrence of population inversion. The laser emission at different pump energies (**Figure 11**) shows that the spectral bandwidth becomes broader when the

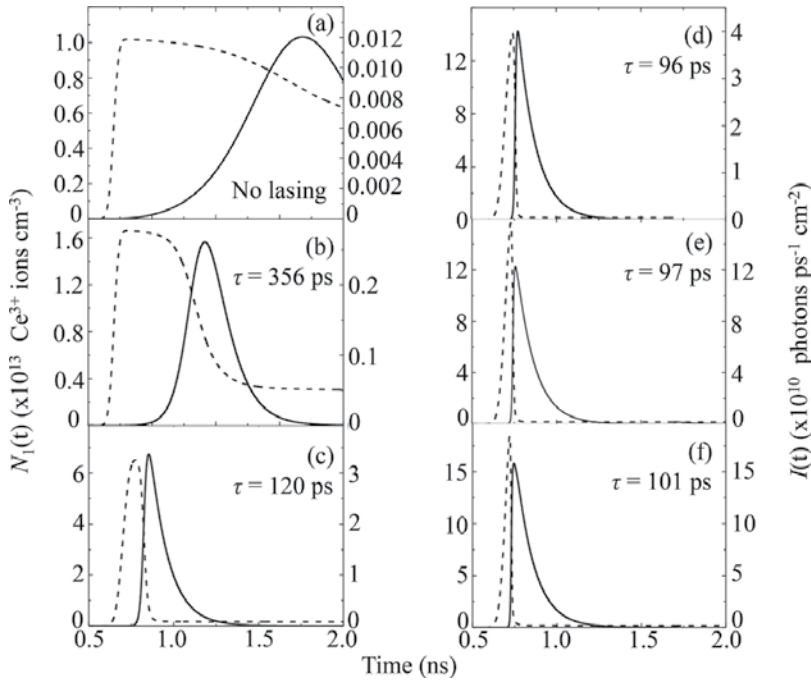


Figure 10. Temporal evolution of the Ce:LiCAF laser emission for different pump energies: (a) 25 μJ , (b) 37.5 μJ , (c) 150 μJ , (d) 360 μJ , (e) 450 μJ , and (f) 600 μJ . The dashed plot represents $N_1(t)$ while the solid plot represents $I(t)$. The following parameters were kept constant: $N = 5 \times 10^{17} \text{ cm}^{-3}$, $l = 10 \text{ mm}$, $L = 15 \text{ mm}$, $R_1 = 100\%$, $R_2 = 30\%$, $r = 100 \mu\text{m}$, $\Delta t = 75 \text{ ps}$, and $n = 1.41$.

pump energy is increased. This is expected since wavelengths close to 288 nm have sufficient energy to achieve considerable gain.

The laser output energy was calculated using

$$E_{out} = \int_0^t \int_{\lambda_1}^{\lambda_n} \frac{(1 - R_2)\pi r^2 hc}{\lambda} I(\lambda, t) \quad (5)$$

where $I(\lambda, t)$ is the laser intensity at wavelength λ and time t . As expected, the output energy increases as the pump energy is increased. This trend is shown in **Figure 12**.

Figures 10 and **12** indicate that in theory, increasing the pump energy would result to an ultrashort (ps pulse duration) laser pulse with micro Joule energy. However, absorption saturation and the damage threshold of the crystal limit the choice of pump energy in experiments. Eq. (4) is therefore integrated in order to determine where absorption saturation of the 266-nm pump begins for the crystal used in experiments. The crystal is 1-cm long and the total number of Ce^{3+} ions is 6.3×10^{14} . The absorption saturation is determined to begin at 1.4 mJ pump energy when the beam spot radius is 100- μm . However, the damage threshold of Ce:LiCAF at 266 nm (wavelength of pump) is about 2 J/cm² [76]. Therefore, the maximum pump energy

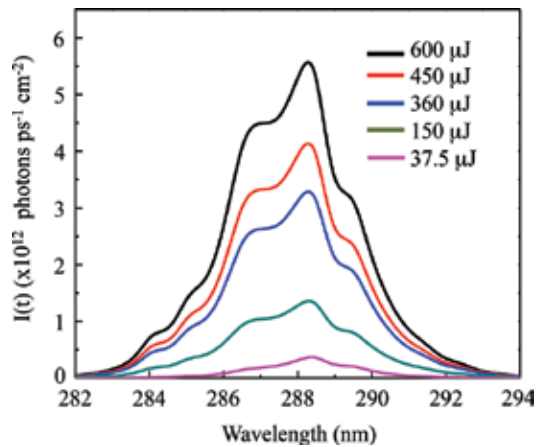


Figure 11. Spectral profiles of the laser emissions for the different pump energies shown in **Figure 10**.

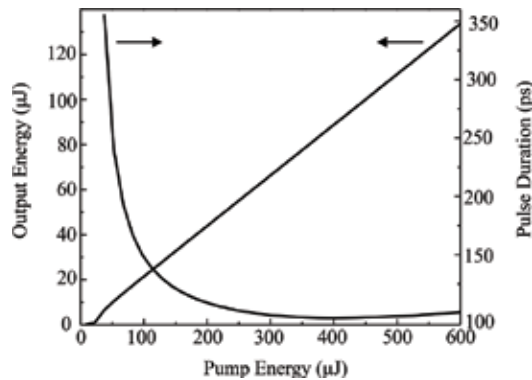


Figure 12. Output energy and pulse duration for different pump energies. Calculation parameters used simulate experimental conditions: $N = 5 \times 10^{17} \text{ cm}^{-3}$, $l = 10 \text{ mm}$, $L = 15 \text{ mm}$, $R_1 = 100\%$, $R_2 = 30\%$, $r = 100 \text{ } \mu\text{m}$, $\Delta t = 75 \text{ ps}$, and $n = 1.41$. The slope efficiency is about 23%.

used in the simulations is $600 \text{ } \mu\text{J}$ as the damage threshold is reached at this amount of energy for a $100 \text{ } \mu\text{m}$ -beam radius. The laser cavity that was used in Ref. [71] amounts to a laser emission efficiency of about 23% as indicated by **Figure 12**. The same figure (**Figure 12**) also shows that a $360 \text{ } \mu\text{J}$ -pump energy could generate 96 ps pulses. This would be the shortest pulse duration. From this observation, the spectro-temporal evolution of the laser pulse emission at $360 \text{ } \mu\text{J}$ -pump energy is calculated and shown in **Figure 13**.

On the other hand, the suppression of the laser emission is evident where the Ce:LiCAF's gain coefficient is lower, particularly on either side of the 288.5-nm wavelength. Because of this, the maximum intensity of the broadband UV laser emission is expected at about 288.5 nm. **Figure 14** confirms that the gain coefficient is still maximum at around 288.5 nm and hence, a shift in the position of the laser peak is not expected. **Figure 10d** shows the temporal evolution of the laser pulse from 284 to 293 nm. This wavelength range spans the whole bandwidth of the UV laser

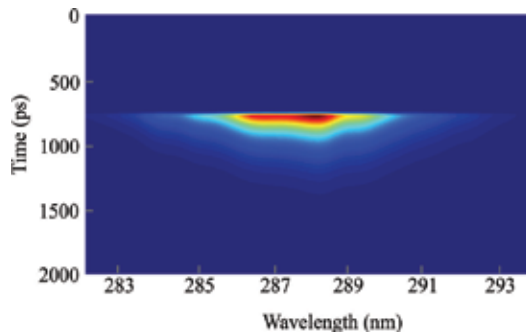


Figure 13. Spectro-temporal evolution of the shortest possible laser pulse duration (96 ps) achievable using experimental resonator parameters $N = 5 \times 10^{17} \text{ cm}^{-3}$, $l = 10 \text{ mm}$, $L = 15 \text{ mm}$, $R_1 = 100\%$, $R_2 = 30\%$, $r = 100 \mu\text{m}$, $\Delta t = 75 \text{ ps}$, $n = 1.41$, and $360\text{-}\mu\text{J}$ pump energy.

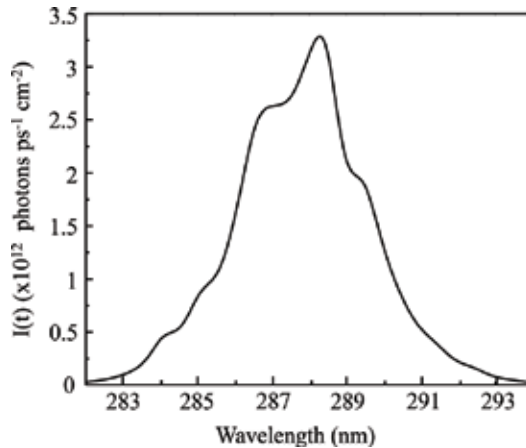


Figure 14. Spectral dynamics of the shortest achievable laser pulse duration (96 ps). The corresponding temporal dynamics is shown in **Figure 10d**.

emission from Ce:LiCAF. By measuring the full width at half maximum (FWHM), the duration of the laser pulse is estimated to be about 96 ps. These calculations indicate that the shortest pulse duration that can be achieved without damaging the crystal is around 96 ps with a slope efficiency of around 23% assuming that the same parameters for the laser cavity and gain medium from the experiment of reference [71] are used.

The design of the laser oscillator cavity is important for optimizing laser emission especially for solid-state gain media. The cavity transient method, which uses the relationship between the cavity round trip time and the fluorescence decay time of the laser gain medium offers a simplified means of generating short laser pulses directly from a cavity that is optically pumped [34, 35]. With this method, factors such as the cavity lifetime of the photon, energy of the pump, and duration of the laser pump pulse strongly influences the pulse duration of laser emission. Moreover, the cavity lifetime of the photon (τ_c) which is described by Eq. (6) is determined by the length of the cavity and the reflectivity of the mirrors.

$$\tau_c = \frac{L + n(l - 1)}{c(1 - \ln(R_1 R_2))} \quad (6)$$

The photon cavity lifetime (τ_c) decreases as the cavity length, L , or the mirror reflectivities, R_1 and R_2 , decrease. Short-pulse laser emission in solid-state gain media can be achieved through the combination of a photon cavity lifetime (τ_c) that is smaller compared to the duration of the pump laser pulse and moderate resonator transients. The latter is brought about by the interaction between the photons in the cavity and the excess population inversion. Previous works have reported using resonator transients in dye lasers to obtain laser pulse durations that are an order of magnitude shorter than the pulse duration of the pump laser [34–36].

In order to extend the technique of resonator transients to solid-state gain media, the rate equations were solved for a variety of output coupler reflectivities and cavity lengths. The optimum condition for setting up a transient cavity was first determined by using a constant value for the output coupler reflectivity while varying the cavity lengths from $L = 2$ mm to $L = 10$ mm. The following summarizes the values of the parameters that were kept constant: $R_1 = 100\%$, $R_2 = 30\%$, $l = 1$ mm, $r = 100$ μm , pumping energy = 140 μJ , and $\Delta t = 75$ ps. The total number of Ce^{3+} ions in the crystal considered here is 6.3×10^{13} . By integrating Eq. (4), the absorption saturation of the 266 nm pump begins at pump energy of 142.8 μJ . Therefore, the maximum energy used in all calculations involving this 1-mm crystal, including **Figures 15–18**, is 140 μJ . Note that the 30% output coupler reflectivity is the same as what was used in the experiment of reference [71]. Calculations show that shorter pulse durations are obtained when the cavity length is shortened. These results are presented in **Figure 15a–d**. It should be noted

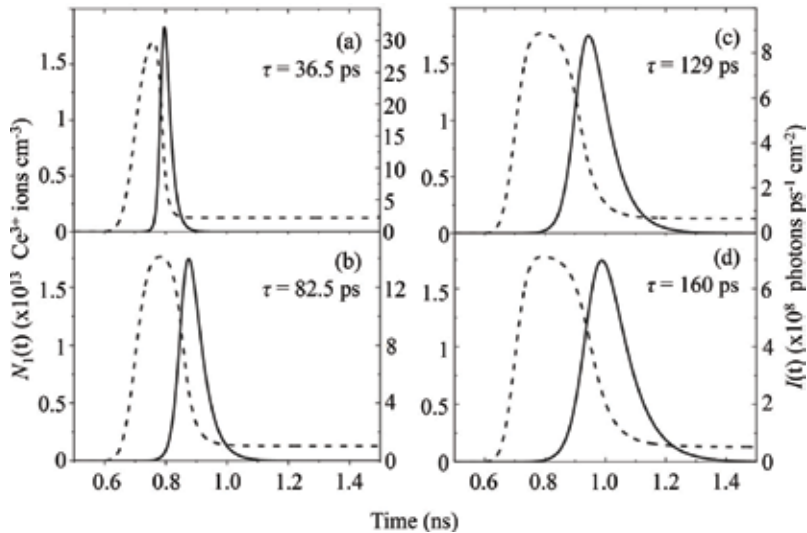


Figure 15. Temporal evolution of the Ce:LiCAF laser emission for different cavity lengths, L , of (a) 2 mm (b) 5 mm, (c) 8 mm, and (d) 10 mm. $I(t)$ is represented by the solid plot while $N_1(t)$ is represented by the dashed plot. The values of the following parameters that were kept constant are: $N = 5 \times 10^{17} \text{ cm}^{-3}$, $l = 1$ mm, $R_1 = 100\%$, $R_2 = 30\%$, $r = 100$ μm , $\Delta t = 75$ ps, $n = 1.41$, and pump energy = 140 μJ . The pump energy was chosen based on the effect of pump energy on pulse duration shown in **Figure 16**.

that the calculation model does not consider energy loss due to cavity length and therefore, the numerical results could over estimate experimental results particularly when the cavity length is long. Regardless of the over estimation, the numerical results show that it is favorable to have a shorter cavity length in order to achieve a shorter pulse durations. As Eq. (5) predicts, smaller photon cavity lifetimes are obtained from smaller cavity lengths and as a consequence, shorter laser pulse durations are also obtained. However, the length of the crystal, l , dictates the limit on the practical size of the cavity, although it would appear from Eq. (5) that ultrashort pulses could be obtained by using ultrashort cavity lengths. From the point of view of crystal growth, 1 mm is a practical crystal length for a LiCAF crystal doped with 1 mol% Ce^{3+} , and for fluoride crystals with 1 mol% rare earth doping in general. A 140- μJ pump energy was used based on **Figure 16**, which shows how the pump energy affects the pulse duration. The output energy is about the

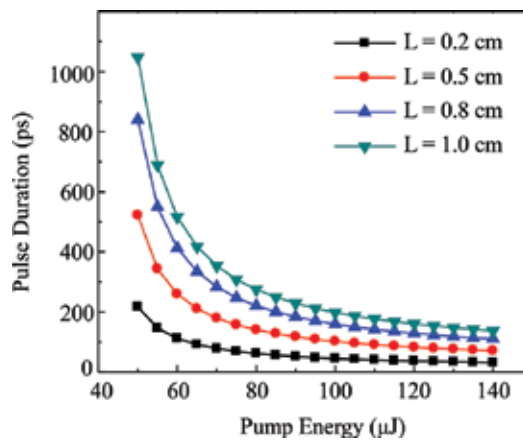


Figure 16. Effect of pump energy on the pulse duration of the Ce:LiCAF laser emission.

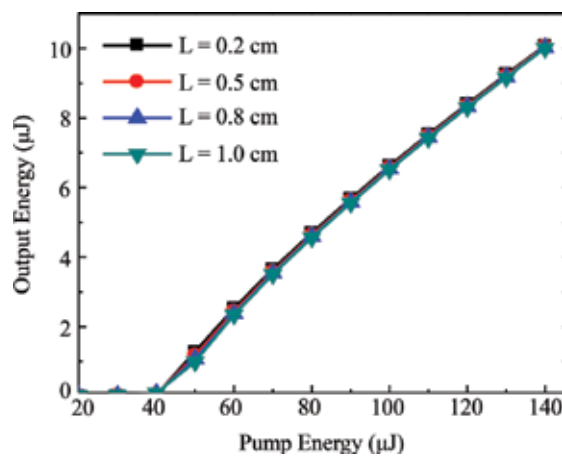


Figure 17. Output energy for the different cavity lengths considered.

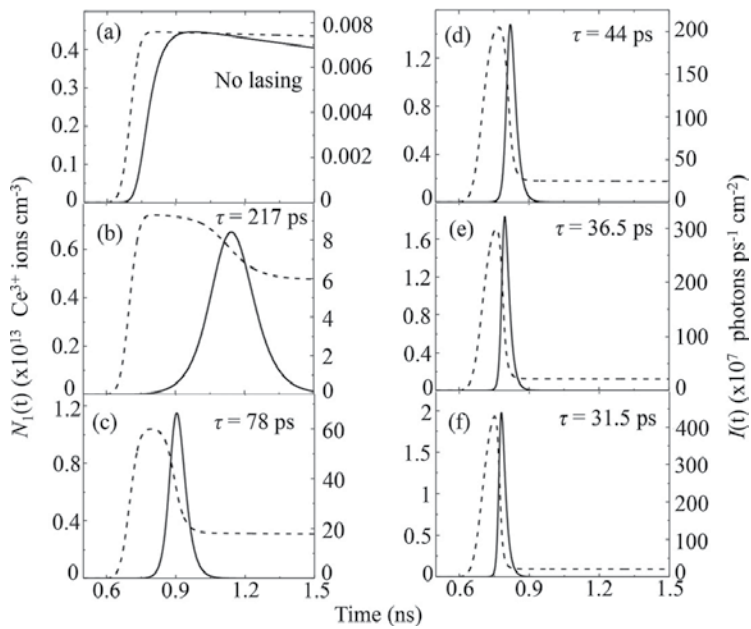


Figure 18. Temporal profile of the laser pulse for different pump energies using a $L = 2$ mm laser oscillator (short cavity). The energies used are: (a) $30 \mu\text{J}$, (b) $50 \mu\text{J}$, (c) $70 \mu\text{J}$, (d) $100 \mu\text{J}$, (e) $120 \mu\text{J}$ and (f) $140 \mu\text{J}$. $I(t)$ is represented by the solid plot while $N_1(t)$ is represented by the dashed plot. The values of the following parameters that were kept constant are: $N = 5 \times 10^{17} \text{ cm}^{-3}$, $l = 1$ mm, $R_1 = 100\%$, $R_2 = 30\%$, $r = 100 \mu\text{m}$, $\Delta t = 75$ ps, and $n = 1.41$. **Figure 19** summarizes the laser output energy and pulse duration for the energies used in **Figure 18**. A duration of 31.5 ps is the shortest pulse duration that can be obtained for $L = 2$ mm and $l = 1$ mm (short resonator cavity). The damage threshold and the absorption saturation of the crystal limit this pulse duration. Therefore, optimizing growth conditions and doping levels can achieve shorter pulse durations. Nevertheless, this is much shorter than the pulse duration that can be obtained using the same values for the pump energy ($140 \mu\text{J}$), and output coupler reflectivity (30%), but with a cavity length of $L = 15$ mm and a crystal length of $l = 10$ mm. As **Figure 12** shows, this laser oscillator that is longer produces pulse duration of 123 ps. On the other hand, a lower output energy is obtained for the same pump energy when the length of the crystal and hence the length of the cavity are shortened. For instance, about $10 \mu\text{J}$ output energy is obtained using $L = 2$ mm and $l = 1$ mm whereas about $30 \mu\text{J}$ output energy is obtained using $L = 15$ mm and $l = 10$ mm for the same $140 \mu\text{J}$ pump energy.

same for the cavity lengths considered ($L = 2\text{--}10$ mm) as shown in **Figure 17**. The slope efficiency is about 10%.

Figure 18a–f shows the temporal evolution of the laser pulse from a short cavity ($L = 2$ mm) for various pump energies. The same crystal and laser cavity parameters used in simulating **Figure 15a–d** were used in **Figure 18a–f**. A single laser pulse with ps pulse duration is achieved when the energy of the pump laser is varied from $50 \mu\text{J}$ to $140 \mu\text{J}$, with a 31.5 ps pulse duration obtained at $140 \mu\text{J}$. Even though it is not shown in the figure, it is worth noting that when the energy of the pump is greater than 3 mJ, lasing occurs almost immediately. This pump energy leads to excess population inversion, as it is much higher than the lasing threshold. As a result of the interaction between the excess population inversion and the photons in the cavity, resonator transients are formed and these are manifested as damped relaxation oscillation or spiking in the laser pulse profile. The laser pulse will eventually approach the shape of the pump pulse when the energy of the pump is increased further,

although the resonator transients (spikes) will still be visible. These spikes have not been observed experimentally since the pump energy in experiments is not high enough.

As discussed earlier, the pump energy and pulse duration of the laser pump as well as the photon cavity lifetime which is determined by the length of the oscillator cavity and the reflectivity of the mirrors strongly influences the pulse duration of the resulting laser emission. Therefore, numerical simulations are performed to quantify the effect of the reflectivity of the output coupler on the laser pulse duration for various pump energies. The results are shown in **Figure 20**. Other parameters are kept constant as follows: $N = 5 \times 10^{17} \text{ cm}^{-3}$, $L = 2 \text{ mm}$, $l = 1 \text{ mm}$, $R_1 = 100\%$, $r = 100 \text{ }\mu\text{m}$, and $\Delta t = 75 \text{ ps}$. It can be observed that a low-Q cavity results to a short laser pulse. The shortest pulse duration is about 31.5 ps when the output coupler

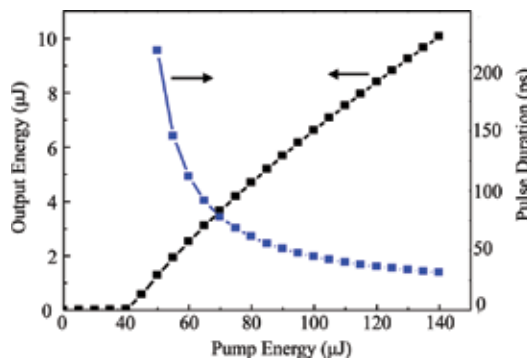


Figure 19. Output energy and pulse duration of the laser pulse from a short cavity oscillator ($L = 2 \text{ mm}$) for different pump energies. The following parameters were kept constant: $N = 5 \times 10^{17} \text{ cm}^{-3}$, $l = 1 \text{ mm}$, $R_1 = 100\%$, $R_2 = 30\%$, $r = 100 \text{ }\mu\text{m}$, $\Delta t = 75 \text{ ps}$, and $n = 1.41$. The slope efficiency is about 24%.

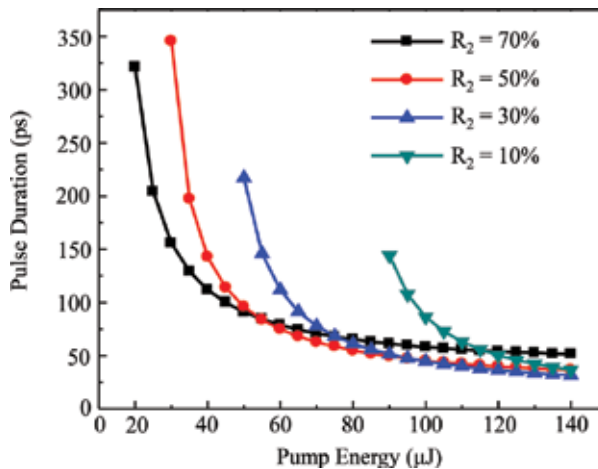


Figure 20. Dependence of pulse duration on pump energy for different output coupler reflectivities: 10, 30, 50, and 70%. A short cavity oscillator ($L = 2 \text{ mm}$) was assumed and the other parameters were kept constant as follows: $N = 5 \times 10^{17} \text{ cm}^{-3}$, $l = 1 \text{ mm}$, $R_1 = 100\%$, $r = 100 \text{ }\mu\text{m}$, $\Delta t = 75 \text{ ps}$, and $n = 1.41$.

reflectivity is 30%. Taking a closer look at the temporal dynamics of the laser pulse for various output coupler reflectivity (**Figure 21a-d**) shows that the threshold for laser emission is reached earlier when the reflectivity is increased. As a result, the laser pulse duration will be longer when the output coupler is highly reflecting. Consequently, a low-Q laser resonator that is established using an output coupler with low reflectivity is required in order to generate short laser pulses low-Q cavity is therefore desirable for generating short-pulse laser emission. **Figure 21a-d** was simulated using laser pump energy of 140 μJ . The choice of pump energy is

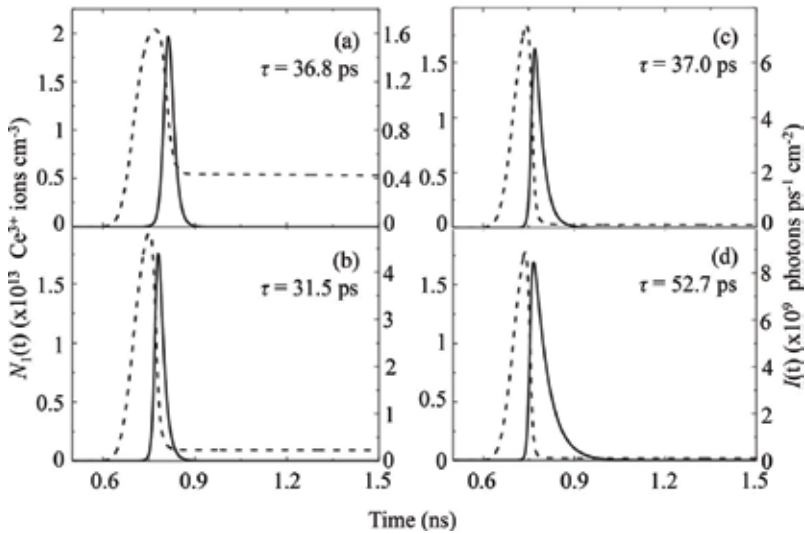


Figure 21. Temporal evolution of the laser emission for various output coupler reflectivity (R_2). The length of the short Ce: LiCAF cavity oscillator is $L = 2$ mm. The various R_2 values considered are: (a) 10%, (b) 30%, (c) 50%, and (d) 70%. The following parameters were kept constant: $N = 5 \times 10^{17} \text{ cm}^{-3}$, $l = 1$ mm, $R_1 = 100\%$, $\Delta t = 75$ ps, $r = 100 \mu\text{m}$, $n = 1.41$, and pumping energy is 140 μJ . The pump energy was chosen based on results in **Figures 19** and **20**. Solid graph represents $I(t)$ while dashed graph represents $N_1(t)$.

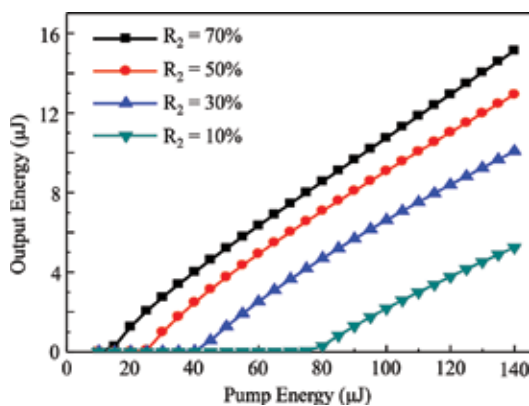


Figure 22. Output energy as a function of pump energy for different output coupler reflectivities ranging from 10 to 70%. A short cavity oscillator ($L = 2$ mm) was assumed and the other parameters were kept constant as follows: $N = 5 \times 10^{17} \text{ cm}^{-3}$, $l = 1$ mm, $R_1 = 100\%$, $r = 100 \mu\text{m}$, $\Delta t = 75$ ps, and $n = 1.41$.

based on the results shown in **Figures 19** and **20**, where pulse duration decreases with increasing energy within the limits of absorption saturation. Theoretically, shorter pulse durations can be achieved using an output coupler with less than 10% reflectivity. Practically, the threshold energy and the slope efficiency limit the choice of reflectivity. As **Figure 22** shows, higher pump energies are needed to achieve lasing in a low-Q laser resonator. If the reflectivity of the output coupler is 10%, for instance, the lasing threshold for obtaining a 31.5-ps laser pulse is 80 μJ and the slope efficiency is 8%. However, increasing the output coupler reflectivity to 30% decreases the threshold energy to 40 μJ and increases the slope efficiency to 10%.

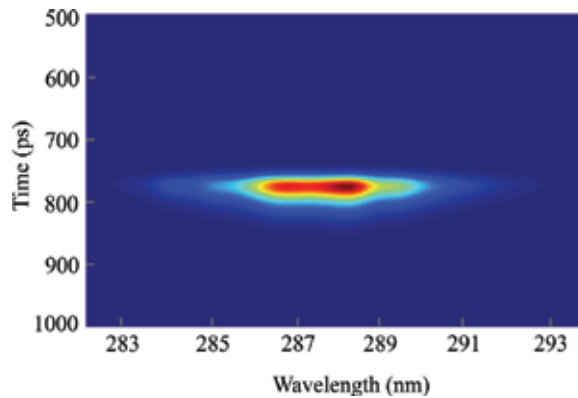


Figure 23. Spectro-temporal evolution of the broadband, short-pulse Ce:LiCAF laser emission from an optimized low-Q ($R_2 = 30\%$), short cavity ($L = 2$ mm) oscillator. A short laser pulse with about 31.5-ps pulse duration, broadband emission centered at 288.5-nm wavelength, and 10 μJ output energy can be obtained practically from a 1-mm long, 1 mol% Ce^{3+} -doped LiCAF crystal when pumped by a 266-nm, 75-ps pump pulse with 140 μJ pump energy. A slope efficiency of about 10% is also feasible with pump energies that are far from the crystal's absorption saturation and damage threshold.

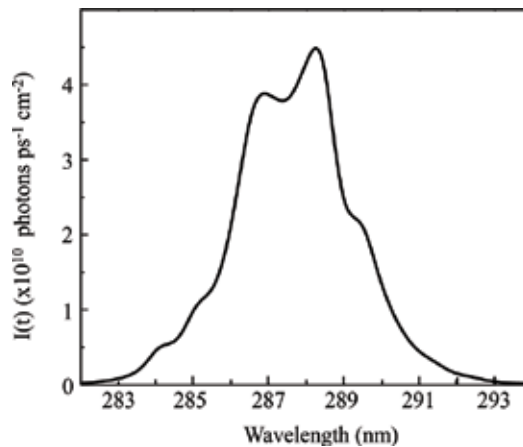


Figure 24. Spectral profile of the broadband, short-pulse Ce:LiCAF laser emission from an optimized low-Q ($R_2 = 30\%$), short cavity ($L = 2$ mm) oscillator. Temporal dynamics is shown in **Figure 21b**.

According to **Figures 19–22**, the optimal transient cavity laser resonator has a 2-mm long cavity and a 30% output coupler reflectivity. These figures also indicate that about 31.5-ps laser pulse duration and about 10% slope efficiency is possible when a 1 mol% Ce-doped LiCAF crystal that is 1 mm long is excited by a 266-nm wavelength pump laser with 75-ps pulse duration and 140 μJ pump energy. These conditions already take into account the crystal's damage threshold, which is about 600- μJ of pump energy for a 100- μm -beam radius as well as its absorption saturation, which is about 142.8 μJ . The 31.5-ps pulse duration is significantly shorter than the experimental pulse duration obtained by reference [71]. The spectro-temporal and the spectral profiles of the broadband 31.5-ps laser pulse are shown in **Figures 23** and **24**, respectively. The spectral profile is consistent with the trend observed in **Figure 11**, regardless of resonator cavity parameters. Maximum gain coefficient is also achieved at around 188.5 nm.

6. Conclusion

In summary, the transient cavity method is extended to a solid-state gain medium. Numerical simulations show that the same principles used to generate ultrashort laser pulses in dye lasers using this technique can be applied to solid-state gain media to generate ultrashort broadband pulses in the UV region. The laser gain medium was represented as a system of two homogeneous broadened singlet states and the numerical simulations solved the laser rate equations for broadband emission. The spectral and temporal evolution of the resulting laser emission was investigated in order to find the optimal cavity length and output coupler reflectivity that will give rise to the formation of resonator transients in the laser oscillator cavity. The calculations reveal that a laser oscillator with a short cavity and a low Q is ideal for the formation of resonator transients, which then lead to ultrashort (ps) laser emission. Specifically, a 2-mm cavity length and a 10% output coupler reflectivity can be used to generate a single 31.5 ps pulse using a 1-mm long Ce:LiCAF crystal with 1 mol% Ce^{3+} ion doping concentration. Although this work used Ce:LiCAF crystal as the laser gain medium, the transient cavity method can also be applied to generate ultrashort laser pulses using other rare earth-doped fluoride crystals.

Acknowledgements

This research was supported by the Massey University Research Fund 2018 (MURF 2018 Project No. 1000020752), Institute of Laser Engineering, Osaka University Collaborative Research Grant (Grant no. 2017B1-RADUBAN), JSPS-VAST Joint Research Project (2011–2014), the Vietnam National Foundation for Science and Technology Development (NAFOSTED) under Grant Numbers 103.06.89.09 and 103.03-2015.29. M. Cadatal-Raduban and M.V. Luong are very grateful to K.G. Steenbergen and P. Schwerdtfeger for their input and valuable discussions on the numerical simulation of the electronic properties of the LiCAF and LiSAF host.

Author details

Marilou Cadatal-Raduban^{1*}, Minh Hong Pham², Luong Viet Mui³, Nguyen Dai Hung² and Nobuhiko Sarukura³

*Address all correspondence to: m.raduban@massey.ac.nz

1 Centre for Theoretical Chemistry and Physics, Institute of Natural and Mathematical Sciences, Massey University, Auckland, New Zealand

2 Institute of Physics, Vietnam Academy of Sciences, Hanoi, Vietnam

3 Institute of Laser Engineering, Osaka University, Suita, Osaka, Japan

References

- [1] Link S, Durr HA, Eberhardt W. Femtosecond spectroscopy. *Journal of Physics: Condensed Matter*. 2001;**13**:7873-7884. PII: S0953-8984(01)26618-1
- [2] Sarukura N, Liu Z, Ohtake H, Segawa Y, Dubinskii MA, Semashko VV, Naumov AK, Korableva SL, Abdulsabirov RY. Ultraviolet short pulses from an all-solid-state Ce:LiCAF master-oscillator power-amplifier system. *Optics Letters*. 1997;**22**:994-996. DOI: 10.1364/OL.22.000994
- [3] Assion A, Baumert T, Bergt M, Brixner T, Kiefer B, Seyfried V, Strehle M, Gerber G. Control of chemical reactions by feedback-optimized phase-shaped femtosecond laser pulses. *Science*. 1998;**282**:919-922. DOI: 10.1126/science.282.5390.919
- [4] Baltuška A, Udem T, Uiberacker M, Hentschel M, Goulielmakis E, Gohle C, Holzwarth R, Yakovlev VS, Scrinzi A, Hänsch TW, Krausz F. Attosecond control of electronic processes by intense light fields. *Nature*. 2003;**421**:611-615. DOI: 10.1038/nature01414
- [5] Nam KB, Li J, Kim KH, Lin JY, Jiang HX. Growth and deep ultraviolet picosecond time-resolved photoluminescence studies of AlN/GaN multiple quantum wells. *Applied Physics Letters*. 2001;**78**:3690-3692. DOI: 10.1063/1.1377317
- [6] Miyamoto I, Horn A, Gottmann J. Local melting of glass material and its application to direct fusion welding by Ps-laser pulses. *Journal of Laser Micro/Nanoengineering*. 2007;**2**:7-14. DOI: 10.2961/jlmn.2007.01.0002
- [7] Pissadakis S, Konstantaki M. Photosensitivity of germanosilicate fibers using 213 nm, picosecond Nd:YAG radiation. *Optics Express*. 2005;**13**:2605-2610. DOI: 10.1364/OPEX.13.002605
- [8] Chen Y, Vertes A. Adjustable fragmentation in laser desorption/ionization from laser-induced silicon microcolumn arrays. *Analytical Chemistry*. 2006;**78**:5835-5844. DOI: 10.1021/ac060405n

- [9] Raciukaitis G, Jacinavicius A, Brikas M, Balickas S. "Picosecond lasers in micromachining," Proceedings of ICALEO, Laser Microfabrication Conference, 2003, p. M308. <http://lia.scitation.org/doi/10.2351/1.2164480>
- [10] Weickhard C, Tonnie K. Short pulse laser mass spectrometry of nitrotoluen: Ionization and fragmentation behavior. *Rapid Communications in Mass Spectrometry*. 2002;**16**:442-446. DOI: 10.1002/rcm.567
- [11] Roy S, Meyer TR, Gord JR. Time-resolved dynamics of resonant and nonresonant broadband picosecond coherent anti-Stokes Raman scattering signals. *Applied Physics Letters*. 2005;**87**:264103. DOI: 10.1063/1.2159576
- [12] Gedvilas M, Raciukaitis G. Investigation of UV picosecond laser ablation of polymers. In: Proc. SPIE 6157, Workshop on Laser Applications in Europe. 2005. p. 61570T. DOI: 10.1117/12.661141
- [13] Wintner E. Numerical evaluation of optical pumpprobe experiments. *Journal of Applied Physics*. 1985;**57**:1533-1537. DOI: 10.1063/1.334467
- [14] Watanabe S, Endoh A, Watanabe M, Sarukura N, Hata K. Multiterawatt excimer-laser system. *Journal of the Optical Society of America B-Optical Physics*. 1989;**6**:1870-1876. DOI: 10.1364/JOSAB.6.001870
- [15] Watanabe S, Endoh A, Watanabe M, Surakura N. Single-shot measurement of subpicosecond KrF pulse width by three-photon fluorescence of the XeF visible transition. *Optics Letters*. 1988;**13**:996-998. DOI: 10.1364/OL.13.000996
- [16] Ehrlich DJ, Moulton PF, Osgood RM Jr. Ultraviolet solid-state Ce:YLF laser at 325 nm. *Optics Letters*. 1979;**4**:184-186. DOI: 10.1364/OL.4.000184
- [17] Ehrlich DJ, Moulton PF, Osgood RM Jr. Optically pumped Ce:LaF₃ laser at 286 nm. *Optics Letters*. 1980;**5**:339-341. DOI: 10.1364/OL.5.000339
- [18] Dubinskii MA, Abdulsabirov RY, Korableva SL, Naumov AK, Semashko VV. 18th International Quantum Electronics Conference, OSA Technical Digest. Vol. 548. Washington, DC: Optical Society of America; 1992
- [19] Dubinskii MA, Abdulsabirov RY, Korableva SL, Naumov AK, Semashko VV. A new active medium for a tunable solid-state UV laser with an excimer pump. *Laser Physics*. 1994;**4**:480
- [20] Dubinskii MA, Semashko VV, Naumov AK, Abdulsabirov RY, Korableva SL. Active medium for all solid-state tunable W laser. *OSA Proceedings Advanced Solid-State Lasers*. 1993;**15**:195-198. DOI: 10.1364/ASSL.1993.LM5
- [21] Dubinskii MA, Semashko VV, Naumov AK, Abdulsabirov RY, Korableva SL. Spectroscopy of a new active medium of a solid-state UV laser with broadband single-pass gain. *Laser Physics*. 1993;**3**:216-217
- [22] Dubinskii MA, Semashko VV, Naumov AK, Abdulsabirov RY, Korableva SL. Ce³⁺-doped colquirite—A new concept of all-solid-state tunable ultraviolet laser. *Journal of Modern Optics*. 1993;**40**:1-5. DOI: 10.1080/09500349314550011

- [23] Marshall CD, Speth JA, Payne SA, Krupke WF, Quarles GJ, Castillo V, Chai BHT. Ultraviolet laser emission properties of Ce³⁺-doped LiSrAlF₆ and LiCaAlF₆. *Journal of the Optical Society of America B*. 1994;**11**:2054-2065. DOI: 10.1364/JOSAB.11.002054
- [24] Liu Z, Kozeki T, Suzuki Y, Sarukura N. Chirped-pulse amplification of ultraviolet femto-second pulses by use of Ce³⁺:LiCaAlF₆ as a broadband, solid-state gain medium. *Optics Letters*. 2002;**26**:301-303. DOI: 10.1364/OL.26.000301
- [25] Liu Z, Kozeki T, Suzuki Y, Sarukura N, Shimamura K, Fukuda T, Hirano M, Hosono H. Ce³⁺:LiCaAlF₆ crystal for high-gain or high-peak-power amplification of ultraviolet femtosecond pulses and new potential ultraviolet gain medium: Ce³⁺:LiSr_{0.8}Ca_{0.2}AlF₆. *IEEE Journal of Selected Topics Quantum Electronics*. 2001;**7**:542-550. DOI: 10.1109/2944.974225
- [26] Sarukura N, Dubinskii MA, Liu Z, Semashko VV, Naumov AK, Korableva SL, Abdulsabirov RY, Edamatsu K, Suzuki Y, Segawa T. Ce³⁺-activated fluoride crystals as prospective active media for widely tunable ultraviolet ultrafast lasers with direct 10-ns pumping. *IEEE Journal of Selected Topics Quantum Electronics*. 2005;**1**:792-794. DOI: 10.1109/2944.473661
- [27] Coutts DW, McGonigle AJS. Cerium-doped fluoride lasers. *IEEE Journal of Quantum Electronics*. 2004;**40**:1430-1440. DOI: 10.1109/JQE.2004.834775
- [28] Pham MH, Cadatal MM, Tatsumi T, Saiki A, Furukawa Y, Nakazato T, Estacio E, Sarukura N, Suyama T, Fukuda K, Kim KJ, Yoshikawa A, Saito F. Laser quality Ce³⁺:LiCaAlF₆ grown by micro-pulling-down method. *Japanese Journal of Applied Physics*. 2008;**47**:5605-5607. DOI: 10.1143/JJAP.47.5605
- [29] Shimamura K, Baldochi SL, Ranieri IM, Sato H, Fujita T, Mazzocchi VL, Parente CBR, Paiva-Santos CO, Santilli CV, Sarukura N, Fukuda T. Crystal growth of Ce-doped and undoped LiCaAlF₆ by the Czochralski technique under CF₄ atmosphere. *Journal of Crystal Growth*. 2001;**223**:383-388. DOI: 10.1016/S0022-0248(01)00593-0
- [30] Sarukura N, Liu Z, Segawa Y, Semashko VV, Naumov AK, Korableva SL, Abdulsabirov RY, Dubinskii MA. Ultraviolet subnanosecond pulse train generation from an all solid state Ce:LiCAF laser. *Applied Physics Letters*. 1995;**67**:602-604. DOI: 10.1063/1.115402
- [31] Spence DE, Kean PN, Sibbett W. 60-fsec pulse generation from a self-mode-locked Ti:sapphire laser. *Optics Letters*. 1991;**16**:42-44. DOI: 10.1364/OL.16.000042
- [32] Alderighi D, Toci G, Vannini M, Parisi D, Tonelli M. Experimental evaluation of the cw lasing threshold for a Ce:LiCaAlF₆ laser. *Optics Express*. 2005;**13**:7256-7264. DOI: 10.1364/OPEX.13.007256
- [33] Granados E, Coutts DW, Spence DJ. Mode-locked deep ultraviolet Ce:LiCAF laser. *Optics Letters*. 2009;**34**:1660-1662. DOI: 10.1364/OL.34.001660
- [34] Roess D. Giant pulse shortening by resonator transients. *Journal of Applied Physics*. 1966;**37**:2004-2006. DOI: 10.1063/1.1708659

- [35] Lin C, Shank CV. Subnanosecond tunable dye laser pulse generation by controlled resonator transients. *Applied Physics Letters*. 1975;**26**:389-391. DOI: 10.1063/1.88188
- [36] Wyatt R. Transient behaviour of pulsed dye lasers. *Applied Physics*. 1980;**21**:353-359. DOI: 10.1007/BF00895927
- [37] Elias LR, Heaps WS, Yen WM. Excitation of UV fluorescence in LaF doped with trivalent cerium and praseodymium. *Physical Review B*. 1973;**8**:4989-4995. DOI: 10.1103/PhysRevB.8.4989
- [38] Yang KH, DeLuca JA. VUV fluorescence of Nd³⁺, Er³⁺, and Tm³⁺-doped trifluorides and tunable coherent sources from 1650 to 2600 Å. *Applied Physics Letters*. 1976;**29**:499-501. DOI: 10.1063/1.89137
- [39] Hamilton DS. Trivalent cerium doped crystals as tunable laser systems: Two bad apples. In: *Tunable Solid State Lasers*. Vol. 47. Berlin: Springer; 1985. pp. 80-90. DOI: 10.1007/978-3-540-39236-1
- [40] Lim KS, Hamilton DS. Optical gain and loss studies in Ce³⁺:YLiF₄. *Journal of the Optical Society of America B*. 1989;**6**:1401-1406. DOI: 10.1364/JOSAB.6.001401
- [41] Owen JF, Dorain PB, Kobayasi T. Excited-state absorption in Eu²⁺:CaF₂ and Ce³⁺:YAG single crystals at 298 and 77 K. *Journal of Applied Physics*. 1981;**52**:1216-1223. DOI: 10.1063/1.329741
- [42] Hamilton DS, Gayen SK, Pogatshnik GJ, Ghen RD. Optical-absorption and photoionization measurements from the excited states of Ce³⁺:Y₃Al₅O₁₂. *Physical Review B*. 1989;**39**:8807-8815. DOI: 10.1103/PhysRevB.39.8807
- [43] Pogatshnik GJ, Hamilton DS. Excited-state photoionization of Ce³⁺ ions in Ce³⁺:CaF₂. *Physical Review B*. 1987;**36**:8251-8257. DOI: 10.1103/PhysRevB.36.8251
- [44] Dubinskii MA, Cefalas AC, Sarantopoulou E, Spyrou SM, Nicolaidis CA, Abdulsabirov RY, Korableva SL, Semashko VV. Efficient LaF₃:Nd³⁺-based vacuum-ultraviolet laser at 172 nm. *Journal of the Optical Society of America B*. 1992;**9**:1148-1150. DOI: 10.1364/JOSAB.9.001148
- [45] Rambaldi P, Moncorge R, Wolf JP, Pedrini C, Gesland JY. Efficient and stable pulsed laser operation of Ce:LiLuF₄ around 308 nm. *Optics Communications*. 1998;**146**:163-166. DOI: 10.1016/S0030-4018(97)00519-1
- [46] McGonigle AJS, Girard S, Coutts DW, Moncorge R. 10 kHz continuously tunable Ce:LiLuF laser. *Electronics Letters*. 1999;**35**:1640-1641. DOI: 10.1049/el:19991112
- [47] Pinto JF, Rosenblatt GH, Esterowitz L, Quarles GJ. Tunable solid-state laser action in Ce³⁺:LiSrAlF. *Electronics Letters*. 1994;**30**:240-241. DOI: 10.1049/el:19940158
- [48] Dorenbos P. 5d-level energies of Ce³⁺ and the crystalline environment. I. Fluoride compounds. *Physical Review B*. 2000;**62**:15640-15649. DOI: 10.1103/PhysRevB.62.15640

- [49] Ono Y, Nakano K, Shimamura K, Fukuda T, Kajitani T. Structural study of colquiriite-type fluorides. *Journal of Crystal Growth*. 2001;**229**:505-509. DOI: 10.1016/S0022-0248(01)01218-0
- [50] Kuze S, Boulay D, Ishizawa N, Kodama N, Yamaga M, Henderson B. Structures of LiCaAlF₆ and LiSrAlF₆ at 120 and 300 K by synchrotron X-ray single-crystal diffraction. *Journal of Solid State Chemistry*. 2004;**177**:3505-3513. DOI: 10.1016/j.jssc.2004.04.039
- [51] Grzechnik A, Dmitriev V, Weber HP, Gesland JY, Smaalen SV. LiSrAlF₆ with the LiBaCrF₆-type structure. *Journal of Physics: Condensed Matter*. 2004;**16**:3005-3013. DOI: 10.1088/0953-8984/16/18/001
- [52] Kresse G, Hafner J. Ab initio molecular dynamics for liquid metals. *Physical Review B*. 1993;**47**:558-561. DOI: 10.1103/PhysRevB.47.558
- [53] Kresse G, Hafner J. Ab initio molecular-dynamics simulation of the liquid-metal—Amorphous-semiconductor transition in germanium. *Physical Review B*. 1994;**49**:14251-14269. DOI: 10.1103/PhysRevB.49.14251
- [54] Kresse G, Furthmuller J. Efficiency of ab-initio total energy calculations for metals and semiconductors using a plane-wave basis set. *Computational Materials Science*. 1996;**6**:15-50. DOI: 10.1016/0927-0256(96)00008-0
- [55] Kresse G, Furthmuller J. Efficient iterative schemes for ab initio total-energy calculations using a plane-wave basis set. *Physical Review B*. 1996;**54**:11169-11186. DOI: 10.1103/PhysRevB.54.11169
- [56] Blochl PE. Projector augmented-wave method. *Physical Review B*. 1994;**50**:17953-17979. DOI: 10.1103/PhysRevB.50.17953
- [57] Kresse G, Joubert D. From ultrasoft pseudopotentials to the projector augmented-wave method. *Physical Review B*. 1999;**59**:1758-1775. DOI: 10.1103/PhysRevB.59.1758
- [58] Perdew J, Burke K, Ernzerhof M. Generalized gradient approximation made simple. *Physical Review Letters*. 1996;**77**:3865-3868. DOI: 10.1103/PhysRevLett.77.3865
- [59] Perdew J, Burke K, Ernzerhof M. Generalized gradient approximation made simple. *Physical Review Letters*. 1997;**78**:1396. DOI: 10.1103/PhysRevLett.78.1396
- [60] Perdew JP, Ernzerhof M, Burke K. Rationale for mixing exact exchange with density functional approximations. *Journal of Chemical Physics*. 1996;**105**:9982-9985. DOI: 10.1063/1.472933
- [61] Paier J, Hirschl R, Marsman M, Kresse G. The Perdew-Burke-Ernzerhof exchange-correlation functional applied to the G2-1 test set using a plane-wave basis set. *Journal of Chemical Physics*. 2005;**122**:234102-234115. DOI: 10.1063/1.1926272
- [62] Luong MV, Empizo MJF, Cadatal-Raduban M, Arita R, Minami Y, Shimizu T, Sarukura N, Azechi H, Pham MH, Nguyen HD, Kawazoe Y, Steenbergen KG, Schwerdtfeger P. First-principles calculations of electronic and optical properties of LiCaAlF₆ and LiSrAlF₆

- crystals as VUV to UV solid-state laser materials. *Optical Materials*. 2017;**65**:15-20. DOI: 10.1016/j.optmat.2016.09.062
- [63] Shimizu T, Luong MV, Cadatal-Raduban M, Empizo MJF, Yamanoi K, Arita R, Minami Y, Sarukura N, Mitsuo N, Azechi H, Pham MH, Nguyen HD, Ichiyanagi K, Nozawa S, Fukaya R, Adachi S, Nakamura KG, Fukuda K, Kawazoe Y, Steenbergen KG, Schwerdtfeger P. *Applied Physics Letters*. 2017;**110**:141902. DOI: 10.1063/1.4979106
- [64] Setyawan W, Curtarolo S. High-throughput electronic band structure calculations: Challenges and tools. *Computational Materials Science*. 2010;**49**:299-312. DOI: 10.1016/j.commatsci.2010.05.010
- [65] Shiran N, Gektin A, Neicheva S, Weber M, Derenzo S, Kirm M, True M, Shpinkov I, Spassky D, Shimamura K, Ichinose N. Energy transfer in pure and Ce-doped LiCaAlF₆ and LiSrAlF₆ crystals. *Nuclear Instruments and Methods in Physics Research A*. 2005; **537**:266-270. DOI: 10.1016/j.nima.2004.08.023
- [66] Bensalah A, Shimamura K, Nakano K, Fujita T, Fukuda T. Growth and characterization of LiSrGaF₆ single crystal. *Journal of Crystal Growth*. 2001;**231**:143-147. DOI: 10.1016/S0022-0248(01)01440-3
- [67] Shimamura K, Sato H, Bensalah A, Machida H, Sarukura N, Fukuda T. Growth of LiCaAlF single crystals with an extended diameter and their optical characterizations. *Journal of Alloys and Compounds*. 2002;**343**:204-210. DOI: 10.1016/S0925-8388(02)00125-1
- [68] Lawson JK, Payne SA. Excited-state absorption of Eu²⁺-doped materials. *Physical Review B*. 1993;**47**:14003-14010. DOI: 10.1103/PhysRevB.47.14003
- [69] Schaffers KI, Keszler DA. Structure of LiSrAlF₆. *Acta Crystallographica*. 1991;**C47**:18-20. DOI: 10.1107/S0108270190006205
- [70] Liu Z, Ohtake H, Sarukura N, Dubinski M, Semashko VV, Naumov AK, Korableva SL, Abdulsabirov RY. Subnanosecond tunable ultraviolet pulse generation from a low-Q, short-cavity Ce:LiCAF laser. *Japanese Journal of Applied Physics*. 1997;**36**:L1384-L1386. DOI: 10.1143/JJAP.36.L1384
- [71] Liu Z, Sarukura N, Dubinski M, Abdulsabirov RY, Korableva SL. All-solid-state subnanosecond tunable ultraviolet laser sources based on Ce³⁺-activated fluoride crystals. *Journal of Nonlinear Optical Physics and Materials*. 1999;**8**:41-54. DOI: 10.1142/S0218863599000047
- [72] Meyer YH, Benoist D'azy O, Martin MM, Br  h  ret E. Spectral evolution and relaxation oscillations in dye lasers. *Optics Communications*. 1986;**60**:64-68. DOI: 10.1016/0030-4018(86)90118-5
- [73] Hung ND, Segawa Y, Long P, Trung DV. Studies of picosecond spectro-temporal selection lasers using different dyes in microcavities with a two-stage-arrangement. *Applied Physics B*. 1997;**65**:19-23. DOI: 10.1007/s003400050243

- [74] Hung ND, Plaza P, Martin M, Meyer YH. Generation of tunable subpicosecond pulses using low-Q dye cavities. *Applied Optics*. 1992;**31**:7046-7058. DOI: 10.1364/AO.31.007046
- [75] Pham MH, Cadatal-Raduban M, Luong MV, Le HH, Yamanoi K, Nakazato T, Shimizu T, Sarukura N, Nguyen HD. Numerical simulation of ultraviolet picosecond Ce:LiCAF laser emission by optimized resonator transients. *Japanese Journal of Applied Physics*. 2014;**53**: 062701. DOI: 10.7567/JJAP.53.062701
- [76] Ono S, Suzuki Y, Kozeki T, Murakami H, Ohtake H, Sarukura N, Sato H, Machida S, Shimamura K, Fukuda T. High-energy, all-solid-state, ultraviolet laser power-amplifier module design and its output-energy scaling principle. *Applied Optics*. 2002;**41**:7556-7560. DOI: 10.1364/AO.41.007556

Exact Finite Differences for Quantum Mechanics

Armando Martínez-Pérez and Gabino Torres-Vega

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.71956>

Abstract

We introduce a finite difference derivative, on a non-uniform partition, with the characteristic that the derivative of the exponential function is the exponential function itself, times a constant, which is similar to what happens in the continuous variable case. Aside from its application to perform numerical computations, this is particularly useful in defining a quantum mechanical discrete momentum operator.

Keywords: exact finite differences derivative, discrete quantum mechanical momentum operator, time operator

1. Introduction

Even though the calculus of finite differences is an interesting subject on its own [1–4] that scheme is mainly used to perform numerical computations with the help of a computer. Finite differences methods give approximate expressions for operators like the derivative or the integral of functions, and it is expected that we get a good approximation when the separation between the points of the partition is small; the smaller it becomes the better.

The momentum operator of Quantum Mechanics, when considering continuous variables, is related to the derivative of functions, but its form, when the variable takes discrete values, is not known yet (an approach is found in Ref. [5]); we need an exact expression for the momentum operator in discrete Quantum Mechanics. Thus, to have an expression for the quantum mechanical momentum operator on a mesh of points, we need an exact expression for the derivative on a mesh of points. In this chapter, we intend to modify the usual finite differences definition of the derivative on a partition to propose an operator that can be used as a momentum operator for discrete Quantum Mechanics.

2. Exact first-order finite differences derivatives of functions

In this section, we intend to introduce a finite differences derivative, which has the same eigenfunction as for the continuous variable case. We start with results valid for any function, but we will concentrate, later in the chapter, on the exponential function because that function is used to perform translations along several directions in the quantum realm. The resulting derivative operator will depend on the point at which it is evaluated as well as on the partition of the interval and on the function of interest. This is the trade-off for having exact finite differences derivatives.

2.1. Backward and forward finite differences derivatives

An exact, backward, finite differences derivative of an absolutely continuous function $g(x)$ (this class of functions is the domain of the momentum operator in Quantum Mechanics), on a partition $P = \{x_1, x_2, \dots, x_N\}$ of N non-uniformly spaced points $\{x_j\}_1^N$, is defined through the requirement that

$$(D_b g)(x_j) := \frac{g(x_j) - g(x_j - \Delta_{j-1})}{\chi_2(j-1)} = g'(x_j), \quad (1)$$

where $\Delta_j = x_{j+1} - x_j$ and the spacing function $\chi_2(j)$, which is a replacement for the usual spacing function Δ_j , is obtained by solving the above equality for $\chi_2(j)$,

$$\chi_2(j-1) := \frac{g(x_j) - g(x_j - \Delta_{j-1})}{g'(x_j)} = \frac{1}{g'(x_j)} \sum_{k=1}^{\infty} \frac{(-1)^{k-1}}{k!} g^{(k)}(x_j) \Delta_{j-1}^k. \quad (2)$$

This is an expression which is valid for points x_j different from the zeroes of $g'(x)$.

A definition for forward finite differences at x_j is

$$(D_f g)(x_j) := \frac{g(x_j + \Delta_j) - g(x_j)}{\chi_1(j)} = g'(x_j), \quad (3)$$

where

$$\chi_1(j) := \frac{g(x_j + \Delta_j) - g(x_j)}{g'(x_j)} = \frac{1}{g'(x_j)} \sum_{k=1}^{\infty} \frac{1}{k!} g^{(k)}(x_j) (\Delta_j)^k, \quad (4)$$

valid for points different from the zeroes of $g'(x)$.

These definitions coincide with the usual finite differences derivative when the function to which they act on is the linear function $g(x) = a_0 + a_1 x$, $a_0, a_1 \in \mathbb{C}$. An exact finite differences derivative of other functions need of more terms than the one found in the usual definition of a finite differences derivative, as can be seen in Eqs. (2) and (4).

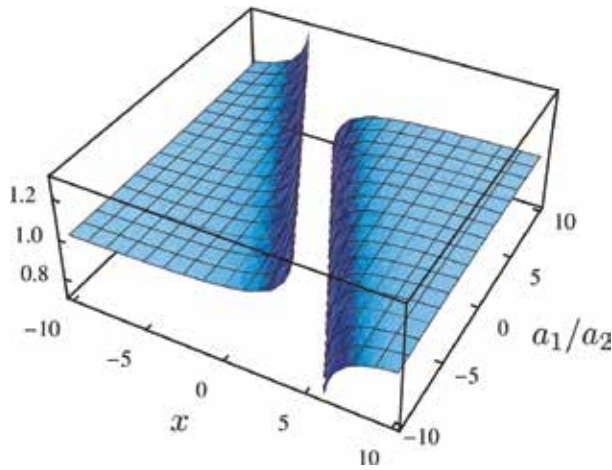


Figure 1. Three-dimensional plot of $\chi_2(x, j)$ for the quadratic function $g(x) = a_0 + a_1x + a_2x^2$ with $\Delta_j = 1$.

Example. For the quadratic function $g(x) = a_0 + a_1x + a_2x^2$, $a_0, a_1, a_2 \in \mathbb{C}$, the spacing function $\chi_2(x, j)$ becomes

$$\chi_2(x, j) = \Delta_j - \frac{\Delta_j^2}{\frac{a_1}{a_2} + 2x}, \tag{5}$$

where $x \neq -a_1/2a_2$. A plot of this function is shown in **Figure 1** for $\Delta_j = 1$.

In the remaining part of this chapter, we only consider the derivative of the exponential function; this choice fixes the form of the spacing functions $\chi_1(j)$ and $\chi_2(j)$.

3. Exact first-order finite differences derivative for the exponential function

Let us consider the exact backward and forward finite differences derivatives of $e^{v x}$, at x_j , given by

$$(D_b e^{v x})_j := \frac{e^{v x_j} - e^{v x_{j-1}}}{\chi_2(v, j-1)} = v e^{v x_j} \text{ and } (D_f e^{v x})_j := \frac{e^{v x_{j+1}} - e^{v x_j}}{\chi_1(v, j)} = v e^{v x_j}, \tag{6}$$

where $v \in \mathbb{C}$ can be a pure real or pure imaginary constant, and the spacing functions $\chi_1(v, j)$ and $\chi_2(v, j)$ are defined as

$$\chi_1(v, j) := \frac{e^{v \Delta_j} - 1}{v} \cong \Delta_j + \frac{v}{2} \Delta_j^2 + O(\Delta_j^3), \tag{7}$$

and

$$\chi_2(v, j) := \frac{1 - e^{-v \Delta_j}}{v} \cong \Delta_j - \frac{v}{2} \Delta_j^2 + O(\Delta_j^3). \tag{8}$$

Note that we recover the usual definitions of a finite differences derivative in the limit $\Delta_j \rightarrow 0$ ($N \rightarrow \infty$) in which case $\chi_1(v, j) = \chi_2(v, j) \rightarrow \Delta_j$. Hereafter, the exact finite differences derivatives that we will consider are

$$(D_b g)_j := \frac{g_j - g_{j-1}}{\chi_2(v, j-1)} \text{ and } (D_f g)_j := \frac{g_{j+1} - g_j}{\chi_1(v, j)}, \tag{9}$$

with $\chi_1(v, j)$ and $\chi_2(v, j)$ given in Eqs. (7) and (8), and some properties of these definitions follow. There is a plot of $\chi_1(v, j)$ in **Figure 2**. The spacing function $\chi_1(v, j)$ is defined for finite values of v and Δ_j .

The summation of a derivative. As is the case for continuous systems, the summation is the inverse operation to the derivative,

$$\sum_{j=n}^m \chi_1(v, j) (D_f g)_j = \sum_{j=n}^m (g_{j+1} - g_j) = g_{m+1} - g_n, \tag{10}$$

where $1 \leq n < m < N$, and

$$\sum_{j=n}^m \chi_2(v, j-1) (D_b g)_j = g_m - g_{n-1}, \tag{11}$$

where $1 < n < m \leq N$.

The exponential function is also an eigenfunction of the summation operation. The usual integral of the exponential function also has its equivalent expression in exact finite differences terms

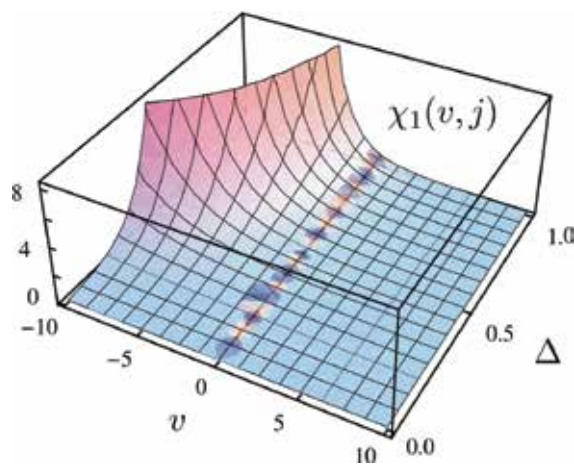


Figure 2. Three-dimensional plot of $\chi_1(v, j)$ for the exponential function $e^{v x}$.

$$\sum_{j=n}^m \chi_1(v, j) v e^{v x_j} = \sum_{j=n}^m \chi_1(v, j) (D_f e^{v x})_j = e^{v x_{m+1}} - e^{v x_n}, \quad (12)$$

where $1 \leq n < m < N$, and

$$\sum_{j=n}^m \chi_2(v, j-1) v e^{v x_j} = e^{v x_m} - e^{v x_{n-1}}, \quad (13)$$

where $1 < n < m \leq N$.

Chain rule. The finite differences versions of the chain rule are

$$\begin{aligned} (D_f g(h(x)))_j &= \frac{g(h(x_{j+1})) - g(h(x_j))}{\chi_1(v, j)} = v \frac{g(h(x_{j+1})) - g(h(x_j))}{e^{v(h(x_{j+1})-h(x_j))} - 1} \frac{e^{v(h(x_{j+1})-h(x_j))} - 1}{v \chi_1(v, j)} \\ &= (D_f g)_j \frac{\chi_1(v, \Delta_{h,j})}{\chi_1(v, j)}, \end{aligned} \quad (14)$$

where

$$(D_f g)_j := v \frac{g(h(x_{j+1})) - g(h(x_j))}{e^{v(h(x_{j+1})-h(x_j))} - 1}, \quad (15)$$

$\chi_1(v, \Delta_{h,j}) = \frac{e^{v(h(x_{j+1})-h(x_j))} - 1}{v}$ and $\Delta_{h,j} = h(x_{j+1}) - h(x_j)$, and

$$(D_b g(h(x)))_j = (D_b g)_j \frac{\chi_2(v, \Delta_{h,j-1})}{\chi_2(v, j-1)}, \quad (16)$$

where $\chi_2(v, \Delta_{h,j-1}) = (1 - e^{-v(h(x_j)-h(x_{j-1}))})/v$,

$$(D_b g)_j := v \frac{g(h(x_j)) - g(h(x_{j-1}))}{1 - e^{-v(h(x_j)-h(x_{j-1}))}}. \quad (17)$$

and $h(x)$ is any absolutely continuous complex function on $[a, b]$.

The derivative of a product of functions. The exact finite differences derivative of a product of functions is

$$\begin{aligned} (D_f gh)_j &= \frac{g_{j+1} h_{j+1} - g_j h_j}{\chi_1(v, j)} = g_{j+1} \frac{h_{j+1} - h_j}{\chi_1(v, j)} + \frac{g_{j+1} - g_j}{\chi_1(v, j)} h_j \\ &= g_{j+1} (D_f h)_j + h_j (D_f g)_j = e^{-v \Delta_j} g_{j+1} (D_b h)_{j+1} + h_j (D_f g)_j, \end{aligned} \quad (18)$$

where $1 \leq j < N$. Also, for the backwards derivative, we have

$$(D_bgh)_j = g_j(D_bh)_j + h_{j-1}(D_bg)_j = g_j(D_bh)_j + e^{v \Delta_j} h_{j-1}(D_fg)_{j-1}, \tag{19}$$

where $1 < j \leq N$.

The derivative of the ratio of two functions. For the finite differences, backward derivative of the ratio of two functions we have

$$\begin{aligned} (D_b \frac{g}{h})_j &= \frac{1}{\chi_2(v, j-1)} \left(\frac{g_j}{h_j} - \frac{g_{j-1}}{h_{j-1}} \right) = \frac{1}{\chi_2(v, j-1)} \left(-\frac{g_j(h_j - h_{j-1})}{h_j h_{j-1}} + \frac{(g_j - g_{j-1})h_j}{h_j h_{j-1}} \right) \\ &= \frac{(D_bg)_j}{h_{j-1}} - g_j \frac{(D_bh)_j}{h_j h_{j-1}}, \end{aligned} \tag{20}$$

$$(D_b \frac{g}{h})_j = \frac{(D_bg)_j}{h_j} - g_{j-1} \frac{(D_bh)_j}{h_j h_{j-1}}, \tag{21}$$

$$(D_f \frac{g}{h})_j = \frac{(D_fg)_j}{h_{j+1}} - g_j \frac{(D_fh)_j}{h_j h_{j+1}}, \tag{22}$$

$$(D_f \frac{g}{h})_j = \frac{(D_fg)_j}{h_j} - g_{j+1} \frac{(D_fh)_j}{h_j h_{j+1}}. \tag{23}$$

Additional properties. A couple of equalities that will be needed below are

$$\frac{1}{\chi_2(v, j)} - \frac{1}{\chi_1(v, j)} = v, \text{ and } \frac{\chi_1(v, j)}{\chi_2(v, j)} = e^{v \Delta_j}. \tag{24}$$

For instance, these equalities imply that

$$(D_fg)_j = e^{-v \Delta_j} (D_bg)_{j+1}. \tag{25}$$

Summation by parts. An important result is the summation by parts. The sum of equalities (18) and (19) combined with equalities (10) and (11) provide the exact finite differences summation by parts results,

$$\sum_{j=n}^m \chi_2(v, j) g_{j+1} (D_bh)_{j+1} + \sum_{j=n}^m \chi_1(v, j) h_j (D_fg)_j = g_{m+1} h_{m+1} - g_n h_n, \tag{26}$$

where $1 \leq j < N$, and

$$\sum_{j=n}^m \chi_2(v, j-1) g_j (D_bh)_j + \sum_{j=n}^m \chi_1(v, j-1) h_{j-1} (D_fg)_{j-1} = g_m h_m - g_{n-1} h_{n-1}, \tag{27}$$

where $1 < j \leq N$.

The integration by parts theorem of continuous functions is the basis that allows to define adjoint, symmetric and self-adjoint operations for continuous variables [8, 9]. Therefore, the

summation by parts results can be used in the finding of an appropriate momentum operator for discrete quantum systems. The summation by parts relates two operators between themselves and with boundary conditions on the functions.

4. The matrix associated to the exact finite differences derivative

It is advantageous to use a matrix to represent the finite differences derivative on the whole interval so that we can consider the whole set of derivatives on the partition at once. Let us consider the backward and forward exact finite differences derivative matrices $\mathbf{D}_{b,f}$ given by

$$\mathbf{D}_b := \begin{pmatrix} \frac{-1}{\chi_1(v,1)} & \frac{1}{\chi_1(v,1)} & 0 & 0 & \dots & 0 & 0 & 0 \\ \frac{-1}{\chi_2(v,1)} & \frac{1}{\chi_2(v,1)} & 0 & 0 & \dots & 0 & 0 & 0 \\ 0 & \frac{-1}{\chi_2(v,2)} & \frac{1}{\chi_2(v,2)} & 0 & \dots & 0 & 0 & 0 \\ 0 & 0 & \frac{-1}{\chi_2(v,3)} & \frac{1}{\chi_2(v,3)} & \dots & 0 & 0 & 0 \\ \vdots & & & & & & & \\ 0 & 0 & 0 & 0 & \dots & \frac{-1}{\chi_2(v,N-2)} & \frac{1}{\chi_2(v,N-2)} & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 & \frac{-1}{\chi_2(v,N-1)} & \frac{1}{\chi_2(v,N-1)} \end{pmatrix} \quad (28)$$

and

$$\mathbf{D}_f := \begin{pmatrix} \frac{-1}{\chi_1(v,1)} & \frac{1}{\chi_1(v,1)} & 0 & 0 & \dots & 0 & 0 & 0 \\ 0 & \frac{-1}{\chi_1(v,2)} & \frac{1}{\chi_1(v,2)} & 0 & \dots & 0 & 0 & 0 \\ 0 & 0 & \frac{-1}{\chi_1(v,3)} & \frac{1}{\chi_1(v,3)} & \dots & 0 & 0 & 0 \\ \vdots & & & & & & & \\ 0 & 0 & 0 & 0 & \dots & \frac{-1}{\chi_1(v,N-2)} & \frac{1}{\chi_1(v,N-2)} & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 & \frac{-1}{\chi_1(v,N-1)} & \frac{1}{\chi_1(v,N-1)} \\ 0 & 0 & 0 & 0 & \dots & 0 & \frac{-1}{\chi_2(v,N-1)} & \frac{1}{\chi_2(v,N-1)} \end{pmatrix} \quad (29)$$

We have used the definition for the backward derivative $(D_b g)_j$ for all the rows of the backward derivative matrix D_b but not for the first line in which we have instead used the forward derivative $(D_f g)_1$. A similar thing was done for the forward derivative matrix D_f . These matrices act on bounded vectors $\mathbf{g} = (g_1, g_2, \dots, g_N)^T \in \mathbb{C}^N$.

The matrix formulation of the derivative operators allows the derivation of some useful results for the derivative itself.

4.1. Higher order derivatives

Many properties can be obtained with the help of the derivative matrices $\mathbf{D}_{b,f}$. Expressions for the exact second finite differences derivative associated to the exponential function are obtained through the square of the derivative matrices $\mathbf{D}_{b,f}$. These expressions are

$$(D_{b,f}^2 g)_1 = v \frac{g_2 - g_1}{\chi_1(v, 1)} = v (D_f g)_1, \tag{30}$$

$$(D_{b,f}^2 g)_2 = v \frac{g_2 - g_1}{\chi_2(v, 1)} = v (D_b g)_2, \tag{31}$$

$$(D_{b,f}^2 g)_j = \frac{1}{\chi_2(v, j - 1)} \left(\frac{g_j - g_{j-1}}{\chi_2(v, j - 1)} - \frac{g_{j-1} - g_{j-2}}{\chi_2(v, j - 2)} \right) = \frac{(D_f g)_j - (D_f g)_{j-1}}{\chi_2(v, j - 1)}, \quad j = 3, \dots, N. \tag{32}$$

These expressions have the exponential function $e^{v \cdot x}$ as one of their eigenfunctions with eigenvalue v^2 , as is also the case of the continuous variable derivative. Higher order derivatives can be obtained in an analogous way.

The derivative matrices are singular, which means that they do not have an inverse matrix, but, at a local level, the inverse operator to the derivative is the summation, as we have already shown in a previous section.

4.2. Eigenfunctions and eigenvectors of $D_{b,f}$

Now that we have the matrices $D_{b,f}$ representing the backward and forward derivatives, we are interested in finding their eigenvalues $\lambda \in \mathbb{C}$ and its corresponding eigenvectors \mathbf{e}_λ . Therefore, we begin by finding the values of λ for which the matrices $D_{b,f} - \lambda \mathbf{I}$ are not invertible, that is, when they are singular.

On one hand, for the backward finite difference matrix D_b , the characteristic polynomial is

$$|D_b - \lambda \mathbf{I}| = \lambda \left(\lambda + \frac{1}{\chi_1(v, 1)} - \frac{1}{\chi_2(v, 1)} \right) \left(\frac{1}{\chi_2(v, 2)} - \lambda \right) \left(\frac{1}{\chi_2(v, 3)} - \lambda \right) \dots \left(\frac{1}{\chi_2(v, N - 1)} - \lambda \right) = 0, \tag{33}$$

whose roots are $\lambda_0 = 0$, $\lambda_v = -1/\chi_1(v, 1) + 1/\chi_2(v, 1) = v$ and $\lambda_j = 1/\chi_2(v, j)$, $2 \leq j \leq N - 1$. Let us denote by $\mathbf{e}_\lambda = (e_{\lambda,1}, e_{\lambda,2}, \dots, e_{\lambda,N})^T$ to the eigenvector corresponding to the eigenvalue λ . The system of equations for the components of the eigenvectors is

$$-\frac{e_{\lambda,k}}{\chi_2(v,k)} + \left(\frac{1}{\chi_2(v,k)} - \lambda \right) e_{\lambda,k+1} = 0, \tag{34}$$

with $k = 1, \dots, N - 1$. Then, the eigenvectors are

$$\mathbf{e}_0 = C \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}, \quad \mathbf{e}_v = C \begin{pmatrix} e^{v \cdot x_1} \\ e^{v \cdot x_2} \\ \vdots \\ e^{v \cdot x_N} \end{pmatrix}, \quad \mathbf{e}_j = C \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ y_{j+1} \\ \vdots \\ y_N \end{pmatrix}, \tag{35}$$

where C is the normalization constant, and

$$y_m = \prod_{k=m}^{N-1} \left(1 - \frac{\chi_2(v, k)}{\chi_2(v, j)} \right), \quad (36)$$

where $2 \leq j < m \leq N$. The quantities y_m usually are very small.

On the other hand, for the forward finite difference matrix, \mathbf{D}_f , the characteristic polynomial is

$$|\mathbf{D}_f - \lambda \mathbf{I}| = (-1)^{N-1} \lambda (v - \lambda) \left(\frac{1}{\chi_1(v, 1)} + \lambda \right) \left(\frac{1}{\chi_1(v, 2)} + \lambda \right) \cdots \left(\frac{1}{\chi_1(v, N-2)} + \lambda \right) = 0, \quad (37)$$

Thus, the eigenvalues for the forward derivative are $\lambda_0 = 0$, $\lambda_v = v$ and $\lambda_j = -1/\chi_1(v, j)$, $1 \leq j \leq N - 2$, and the corresponding eigenvectors are

$$e_0 = C \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}, \quad e_v = C \begin{pmatrix} e^{v x_1} \\ e^{v x_2} \\ \vdots \\ e^{v x_N} \end{pmatrix}, \quad e_j = C \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_j \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \quad (38)$$

where

$$w_m = \prod_{k=1}^{m-1} \left(1 - \frac{\chi_1(v, k)}{\chi_1(v, j)} \right) \quad (39)$$

and $1 \leq m < j \leq N - 2$. The quantities w_m are also very small; in fact, they vanish for the equally spaced partition.

The matrices, $\mathbf{D}_{b,f}$, have the same eigenvalues λ_0 and λ_v , and eigenvectors which are the discretization of the function $g_v(x) = C e^{v x}$ on the partition $\{x_1, x_2, \dots, x_N\}$ (the eigenvector $(1, 1, \dots, 1)^T$ correspond to the eigenvalue $v=0$). This is the same eigenfunction that is found in the continuous variable case because the exponential function is indeed an eigenfunction of the continuous derivative. We note that the local derivatives $(D_{b,f}g)_j$ have the same eigenfunctions as the matrices $\mathbf{D}_{b,f}$ which are global objects. The other eigenvectors are fluctuations around the null vector, which is the trivial eigenvector of the derivative.

4.3. The commutator between coordinate and derivative

Since the following equality holds:

$$(D_b x)_j = \frac{x_j - x_{j-1}}{\chi_2(v, j-1)} = \frac{\Delta_{j-1}}{\chi_2(v, j-1)} = 1 + \frac{v}{2} \Delta_{j-1} + \mathcal{O}(\Delta_{j-1}^2), \quad (40)$$

from a local point of view, we have

$$(D_b xg)_j = x_j(D_b g)_j + g_{j-1}(D_b x)_j = x_j(D_b g)_j + g_{j-1} \frac{\Delta_{j-1}}{\chi_2(v, j-1)}, \tag{41}$$

and then, the commutator between x and D_b , acting on g , is given by

$$([D_b, x]g)_j = g_{j-1} \left(1 + \frac{v}{2} \Delta_{j-1} + \mathcal{O}(\Delta_{j-1}^2) \right).$$

Thus, the commutator between D_b and x becomes one in the limit of small Δ_j , or large N , because it also happens that $g_{j-1} \rightarrow g_j$.

We now consider the commutator between the coordinate matrix $\mathbf{Q} := \text{diag}(x_1, x_2, \dots, x_N)$, and the forward derivative matrix \mathbf{D}_f . That commutator is

$$[\mathbf{D}_f, \mathbf{Q}] = \begin{pmatrix} 0 & \frac{\Delta_1}{\chi_1(v,1)} & 0 & \dots & 0 & 0 & 0 \\ 0 & 0 & \frac{\Delta_2}{\chi_1(v,2)} & \dots & 0 & 0 & 0 \\ \vdots & & & & & & \\ 0 & 0 & 0 & \dots & 0 & \frac{\Delta_{N-2}}{\chi_1(v, N-2)} & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 & \frac{\Delta_{N-1}}{\chi_1(v, N-1)} \\ 0 & 0 & 0 & \dots & 0 & \frac{\Delta_{N-1}}{\chi_2(v, N-1)} & 0 \end{pmatrix} \tag{42}$$

The small Δ_j ($N \rightarrow \infty$) approximation of this commutator is just

$$[\mathbf{D}_f, \mathbf{Q}] \xrightarrow{\Delta_j \rightarrow 0} \begin{pmatrix} 0 & 1 & 0 & \dots & 0 & 0 & 0 \\ 0 & 0 & 1 & \dots & 0 & 0 & 0 \\ \vdots & & & & & & \\ 0 & 0 & 0 & \dots & 0 & 1 & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 & 1 \\ 0 & 0 & 0 & \dots & 0 & 1 & 0 \end{pmatrix} \tag{43}$$

There is coincidence with the local calculation; as expected, this matrix approaches the identity matrix in the small Δ_j limit. Note that this matrix is composed of *backward translations* of the first $N - 1$ points and a *forward translation* of the point $N - 1$ without periodicity; the value of the first point is lost.

4.4. Translations

It is well known that the derivative is the generator of the translations of its domain [8]. Therefore, here we investigate briefly how translations are carried out by means of the derivative matrices $\mathbf{D}_{b,f}$ used as their generators. We will focus on the translation of the common eigenvector $e_v = (e^{v \cdot x_1}, e^{v \cdot x_2}, \dots, e^{v \cdot x_N})^T$ of both matrices.

Let the linear transformation represented by the matrix formed by means of the standard definition of a translation operator and of the exponential operator, given by

$$e^s D_{b,f} = \sum_{k=0}^{\infty} \frac{(s D_{b,f})^k}{k!}, \tag{44}$$

where $s \in \mathbb{R}$. Since e_v is an eigenvector of $D_{b,f}$ with eigenvalue v (see Eqs. (35) and (38)), it follows that $D_{b,f}^k e_v = v^k e_v$, $k = 1, 2, \dots$, and then,

$$e^s D_{b,f} e_v = \sum_{k=0}^{\infty} \frac{(s v)^k}{k!} e_v = e^{s v} e_v = C \begin{pmatrix} e^{v (s+x_1)} \\ e^{v (s+x_2)} \\ \vdots \\ e^{v (s+x_N)} \end{pmatrix}, \tag{45}$$

that is, $e^{s v}$ is an eigenvalue of $e^s D_{b,f}$ with corresponding eigenvector e_v , but the right-hand side of this equality is also a translation by the amount s of the domain of the derivative operators. We point out that s is arbitrary and then the vector $e'_v = e^s D_{b,f} e_v$ is the function $e^{v \cdot x}$ evaluated at the points of the translated partition $P' = \{x_1 + s, x_2 + s, \dots, x_N + s\}$. Thus, we can perform not only discrete translations but continuous translations as well.

The usual periodic, discrete translation found in the papers of other authors [6, 7] is obtained when the separation between the partition points is the same (denoted by Δ , a constant) and with periodic boundary conditions $e_{v,1} = e_{v,N}$.

4.5. Fourier transforms between coordinate and derivative representations

In this section, we define continuous and discrete Fourier transforms and establish some of their properties regarding the Fourier transform of continuous and discrete derivatives. The derivative eigenvalue $-ip$ should be understood, and we will omit it from the formulae below for the sake of simplicity of notation.

Given a function $g(p)$ in the L^1 -space and a non-uniform partition $P = \{x_{-N}, x_{-N+1}, \dots, x_N\}$, with $x_{-N} = -x_N$, the function.

$$(Fg)(x_j) := \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-ix_j p} g(p) dp, \tag{46}$$

is the continuous Fourier transform of $g(p)$ at x_j . Having introduced the summation with weights $\chi_1(v, j)$ of Eq. (10), here we define two discrete Fourier transforms at p as

$$(\mathcal{F}_b g)(p) := \frac{p}{2 \sin(p x_N)} \sum_{j=-N}^N \chi_2(\Delta_j) e^{-ix_{j+1} p} g_{j+1}. \tag{47}$$

$$(\mathcal{F}_f g)(p) := \frac{p}{2 \sin(p x_N)} \sum_{j=-N}^N \chi_1(\Delta_j) e^{-ix_j p} g_j \tag{48}$$

Now, the discrete derivative of the product $e^{-ixp}g$ at x_j , with derivative eigenvalue $-ip$, is readily computed to give (see Eq. (18)).

$$(D_f e^{-ixp}g)_j = -ip g_{j+1}(e^{-ixp})_j + e^{-ixp}(D_f g)_j. \quad (49)$$

The summation of this equality, with weights $\chi_1(-ip, j)$, results in

$$\sum_{j=-N}^{N-1} \chi_1(j)(D_f e^{-ixp}g)_j = -ip \sum_{j=-N}^{N-1} \chi_1(j)g_{j+1}e^{-ixp} + \sum_{j=-N}^{N-1} \chi_1(j)e^{-ixp}(D_f g)_j, \quad (50)$$

or

$$\sum_{j=-N}^{N-1} \chi_1(j)e^{-ixp}(-iD_f g)_j = p \sum_{j=-N}^{N-1} \chi_2(j)e^{-ix_{j+1}p}g_{j+1} + \sum_{j=-N}^{N-1} \chi_1(j)(-iD_f e^{-ixp}g)_j. \quad (51)$$

According to Eqs. (10), (24), (47) and (48), this equality can be rewritten in terms of discrete Fourier transforms.

$$(\mathcal{F}_f(-iD_f g))(p) = p (\mathcal{F}_b g)(p) - ip \frac{e^{-ix_N p}g_N - e^{-ix_{-N} p}g_{-N}}{2 \sin(px_N)}. \quad (52)$$

Another expression for the finite differences of the derivative of a function is obtained as follows. Considering the relationship (see Eq. (18), the second expression with $g = e^{-ixp}$)

$$(D_f e^{-ixp}g)_j = e^{i\Delta_j p} e^{-ix_{j+1}p}(D_b g)_{j+1} + g_j(D_f e^{-ixp})_j. \quad (53)$$

The summation of this equality, with weights $\chi_1(j)$, results in

$$\sum_{j=-N}^{N-1} \chi_2(j)e^{-ix_{j+1}p}(-iD_b g)_{j+1} = p \sum_{j=-N}^{N-1} \chi_1(j)e^{-ix_j p}g_j + \sum_{j=-N}^{N-1} \chi_1(j)(-iD_f e^{-ixp}g)_j, \quad (54)$$

and, according to Eq. (10), this equality can be rewritten as the discrete Fourier transform

$$(\mathcal{F}_b(-iD_b g))(p) = p (\mathcal{F}_f g)(p) - ip \frac{e^{-ix_N p}g_N - e^{-ix_{-N} p}g_{-N}}{2 \sin(px_N)}. \quad (55)$$

These are the equivalent to the well-known identities found in continuous variables theory. Thus, the multiplication by p in forward p -space corresponds to the backward finite differences derivative in coordinate space. Additionally, the multiplication by p in backward p -space corresponds to the forward finite differences derivative in coordinate space, when choosing vanishing or periodic boundary conditions.

The integration by parts of the simple relationship

$$\frac{d e^{ix_j p}}{dp} = i x_j e^{ix_j p} \tag{56}$$

results in

$$i x_j \int_{-\infty}^{\infty} dp e^{ix_j p} h(p) = - \int_{-\infty}^{\infty} dp e^{ix_j p} \frac{d h(p)}{dp} + e^{ix_j p} h(p) \Big|_{p=-\infty}^{\infty} \tag{57}$$

or in terms of continuous Fourier transforms,

$$x_j (Fh)_j = \left(F i \frac{d h}{dp} \right)_j - \frac{i h}{\sqrt{2\pi}} e^{ix_j p} h(p) \Big|_{p=-\infty}^{\infty}. \tag{58}$$

These equalities are like the usual properties between the spaces related by the Fourier transform.

5. Quantum mechanical momentum and time operators

We can apply the results of previous sections to discrete Quantum Mechanics theory. Let us rewrite Eq. (26) in terms of complex wave functions $\psi, \phi \in \ell^2(P, [a, b])$ defined on the partition $P = \{x_1, x_2, \dots, x_N\}$ of $[a, b]$. We obtain

$$\sum_{j=1}^{N-1} \chi_2(v, j) \psi_{j+1}^* (-i\hbar D_b \phi)_{j+1} - \sum_{j=1}^{N-1} \chi_1^*(v, j) \phi_j (-i\hbar D_f \psi)_j = -i\hbar (\psi_N^* \phi_N - \psi_1^* \phi_1). \tag{59}$$

This equality is rewritten as

$$\left(\psi | \widehat{P}_b \phi \right)_b - \left(\widehat{P}_f \psi | \phi \right)_f = -i\hbar (\psi_N^* \phi_N - \psi_1^* \phi_1), \tag{60}$$

where the momentum-like operators \widehat{P}_b and \widehat{P}_f are defined as

$$\widehat{P}_b := -i\hbar D_b, \quad \widehat{P}_f := -i\hbar D_f, \tag{61}$$

and the bilinear forms $(\psi | \phi)_b$ and $(\psi | \phi)_f$ are defined as

$$(\psi | \phi)_b := \sum_{j=1}^{N-1} \chi_2(v, j) \psi_{j+1}^* \phi_{j+1}, \tag{62}$$

$$(\psi | \phi)_f := \sum_{j=1}^{N-1} \chi_1^*(v, j) \psi_j^* \phi_j. \tag{63}$$

We recognize Eqs. (60) and (61) as the finite differences versions of the equation that is used to define the adjoint operator and the symmetry of an operator in continuous Quantum

Mechanics. Thus, we propose that the momentum-like operators \widehat{P}_b and \widehat{P}_f are the “adjoint” of each other, on a finite interval $[a, b]$, when

$$\left(\psi|\widehat{P}_b\phi\right)_b = \left(\widehat{P}_f\psi|\phi\right)_f, \quad (64)$$

together with the boundary condition on the wave functions ψ and ϕ ,

$$\psi_N = e^{i\theta}\psi_1, \quad \phi_N = e^{i\theta}\phi_1, \quad (65)$$

where $\theta \in [0, 2\pi)$ is an arbitrary phase. This gets rid of boundary terms.

With these definitions, we are closer to have a finite differences version of a self-adjoint momentum operator on an interval [12, 13] for use in discrete Quantum Mechanics. We believe that our results will lead to a sound definition of a discrete momentum operator and to the finding of a time operator in Quantum Mechanics [10–13].

6. The particle in a linear potential

As an application of the ideas presented in this chapter, we consider the particle under the influence of the linear potential

$$V(x) = \begin{cases} \infty, & x \leq 0, \\ cx, & x > 0, \end{cases} \quad (66)$$

where $c > 0$. The eigenfunction corresponding to this potential is

$$\psi_E(x) = d \operatorname{Ai} \left[\sqrt[3]{\frac{2mc}{\hbar^2}} \left(x - \frac{E}{c} \right) \right], \quad (67)$$

where Ai denotes the Airy function and d is the normalization factor, m is the mass of the quantum particle and \hbar is Planck’s constant divided by 2π . The boundary condition $\psi_E(x=0) = 0$ provides an expression for the energy eigenvalues E , which is

$$E_n = -\sqrt[3]{\frac{\hbar^2 c^2}{2m}} \alpha_{n+1}, \quad n = 0, 1, \dots \quad (68)$$

where $\{\alpha_n\}$ are the roots of the Airy function, which are negative quantities.

In this case, the energy values are discrete and non-uniformly spaced, and the operator conjugate to the Hamiltonian would be a time-type operator with a discrete derivative $\widehat{T} = -i\hbar D_{b,f}$. The eigenfunctions of this time-type operator are calculated as in Eq. (38)

$$\langle x|t \rangle = \sum_{n=0}^M e^{-i t E_n} \psi_n(x), \quad (69)$$

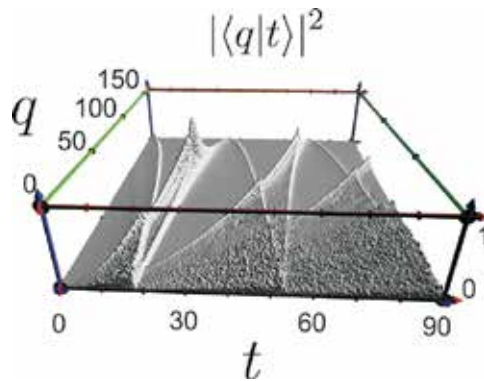


Figure 3. Three-dimensional plot of the squared modulus $|\langle q|t\rangle|^2$ of the time eigenstates for the wall-linear potential in coordinate representation, using 600 energy eigenfunctions. Dimensionless units.

where $\psi_n(x)$ is the eigenfunction of the Hamiltonian with energy E_n , Eq. (67). A plot of these time-type eigenstates, with $M = 600$, is found in **Figure 3**. We can identify the classical trajectories with initial conditions $(x_0, p_0) = (\frac{E_n}{c}, 0)$ in that figure; they are the regions in which the probability is higher. We can also identify the interference pattern between them.

In conclusion, we can have an exact derivative without the need of many terms, and this allows for the definition of adjoint operators related to the derivative on a mesh of points.

Author details

Armando Martínez-Pérez and Gabino Torres-Vega*

*Address all correspondence to: gabino@fis.cinvestav.mx

Physics Department, Cinvestav Mexico City, México

References

- [1] Boole G. A Treatise on the Calculus of Finite Differences. New York: Cambridge University Press; 2009 (1860)
- [2] Harmuth HF, Meert B. Dogma of the Continuum and the Calculus of Finite Differences in Quantum Physics. San Diego: Elsevier Academic Press; 2005
- [3] Jordan C. Calculus of Finite Differences. 2nd ed. New York: Chelsea Publishing Company; 1950
- [4] Richardson CH. An Introduction to the Calculus of Finite Differences. Toronto: D. Van Nostrand; 1954

- [5] Santhanam TS, Tekumalla AR. Quantum mechanics in finite dimensions. *Foundations of Physics*. 1976;**6**:583-587. DOI: 10.1007/BF00715110
- [6] de la Torre AC, Goyeneche D. Quantum mechanics in finite-dimensional Hilbert space. *American Journal of Physics*. 2003;**71**:49-54. DOI: 10.1119/1.1514208
- [7] Pauli W. *General Principles of Quantum Mechanics*. Berlin: Springer-Verlag; 1980
- [8] Martínez-Pérez A, Torres-Vega G. Translations in quantum mechanics revisited. The point spectrum case. *Canadian Journal of Physics*. 2016;**94**:1365-1368. DOI: 10.1139/cjp-2016-0373
- [9] Pauli W. *Handbuch der Physik*. Vol. 23. Berlin: Springer-Verlag; 1926
- [10] Pauli W. *General Principles of Quantum Mechanics*. Berlin: Springer-Verlag; 1980
- [11] Martínez-Pérez A, Torres-Vega G. Translations in quantum mechanics revisited. The point spectrum case. *Canadian Journal of Physics*. 2016;**94**:1365-1368. DOI: 10.1139/cjp-2016-0373
- [12] Gitman DM, Tyutin IV, Voronov BL. *Self-Adjoint Extensions in Quantum Mechanics. General Theory and Applications to Schrödinger and Dirac Equations with Singular Potentials*. New York: Birkhäuser; 2007
- [13] Schmüdgen K. *Unbounded Self-Adjoint Operators on Hilbert Space*. Heidelberg: Springer; 2012

Twin-Grating Fiber Optic Sensors Applied on Wavelength-Division Multiplexing and Its Numerical Resolution

José Trinidad Guillen Bonilla, Héctor Guillen Bonilla,
Antonio Casillas Zamora,
Gustavo Adolfo Vega Gómez,
Nancy Elizabeth Franco Rodríguez,
Alex Guillen Bonilla and Juan Reyes Gómez

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.75586>

Abstract

In this work, the twin-grating fiber optic sensor has been applied on wavelength-division multiplexing. A quasi-distributed sensor formed by three local twin-grating sensors, is numerically simulated. The wavelength channels were 1531.5, 1535.5, and 1539.5 nm. The numerical simulation shows the resolution vs. signal-to-noise rate. Three local twin-grating sensors have approximately the same resolution because all local sensors have the same cavity length and the wavelength channels are very close. All local sensors have two numerical resolutions because the Fourier domain phase analysis algorithm makes two evaluations of the Bragg wavelength shift. The transition between both resolutions can be calculated with the parameters: cavity length, Bragg wavelength channel, refraction index, and enveloped resolution. This transition depends on the noise system, demodulation algorithm, instrumentation, and local sensor properties. A very important point is, a theoretical analysis will permit to know the exact resolution for each local twin-grating sensor.

Keywords: twin-grating fiber optic sensor, wavelength-division multiplexing, numerical resolution, quasi-distributed sensor, numerical simulation

1. Introduction

Optic fiber sensors (OFSs) exhibit small dimensions; they are light weight and made of a dielectric material, vitreous silica. Some measurable parameters are temperature, strain, humidity, pressure, salinity, current, voltage, and concentration. Fiber sensors have as good resolution and accuracy as electronic and mechanical sensors. For this reason, OFSs are very active worldwide. An optic fiber sensor can be extrinsic or intrinsic. In an extrinsic sensor, the fiber acts as a means of getting the light to the sensing localization. In an intrinsic sensor, perturbations act on the fiber and the fiber in turn changes some characteristics of the light inside the fiber [1]. Both sensors types find potential industrial applications. On the other hand, a fiber sensor can also be spatially classified as a distributed sensor, a quasi-distributed sensor, or a point sensor. A distributed sensor is sensitive along its entire length. A quasi-distributed fiber optic sensor is not sensitive along its entire length, but is locally sensitized at various points. A point sensor is sensitive at a specific point along its entire length. In particular, a quasi-distributed sensor uses multiplexing techniques and their combinations. Two fundamental techniques are wavelength-division multiplexing (WDM) and frequency-division multiplexing (FDM). In Ref. [2], Grattan and Sun described the WDM technique:

- The WDM technique received little attention due to the initial high cost of components such as wavelength selective couplers and filters. However, the widespread use of Bragg grating systems has opened up a range of possibilities for the use of wavelength-division multiplexing. **Figure 1a** illustrates a scheme of a quasi-distributed sensor based on the Bragg gratings; its configuration is serial and each Bragg grating has its own Bragg wavelength.

The frequency-division multiplexing scheme [3] is illustrated in **Figure 1b** for a quasi-distributed sensor based on the twin-grating fiber optic sensor. Each twin-grating sensor [4, 5] consists of two identical Bragg gratings and acts as a local sensor. In this configuration, there are m -twin-grating sensors in serial connection. Each interferometer has its own cavity length. However, all Bragg gratings have the same Bragg wavelength to eliminate wavelength-division multiplexing (WDM). The cross-talk noise is eliminated because all Bragg grating had low reflectivity, $r < 1\%$ [5]. In this sensor, the reflection spectrum is the superposition of all frequency components which are produced by all local interferometers. The detection method is known as direct spectrometric detection technique.

Nowadays, the quasi-distributed sensor finds potential application in civil engineering (strain and temperature measurements), industrial process (temperature, strain, level, and pressure measurements), military application (vibration detection), sport science (vibration and strain), and aircraft (strain, vibration, and pressure measurements) [6–9]. This sensor type reduces the cost by sensing point. In this work, a quasi-distributed sensor based on wavelength-division multiplexing and twin-grating sensor is discussed and simulated. The results show the numerical resolution in terms of Bragg wavelength shift. The results demonstrate that twin-grating sensors' resolution is high and the resolution depends on the cavity length.

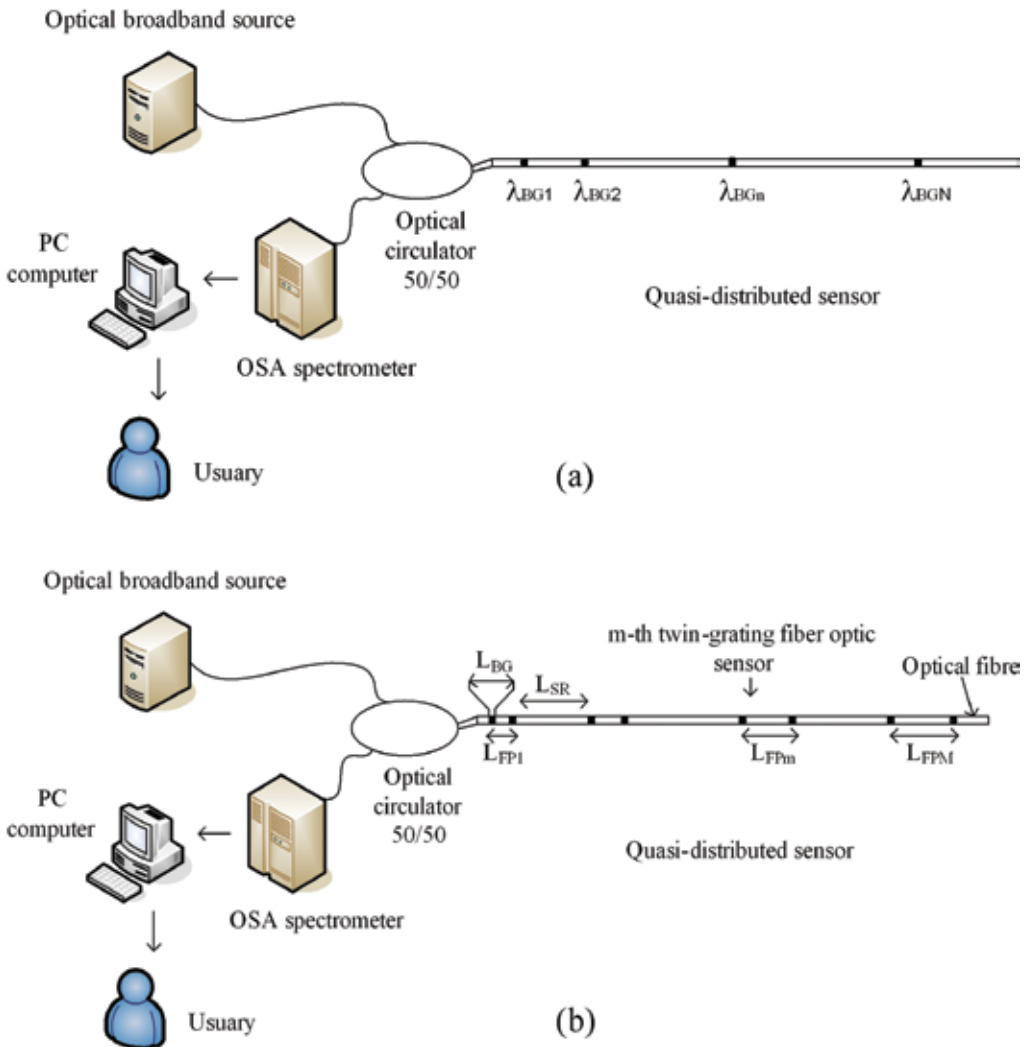


Figure 1. A scheme of a quasi-distributed fiber optic sensor: (a) Bragg gratings and (b) twin-grating interferometers [3].

2. A quasi-distributed fiber optic sensor

Figure 2 illustrates the optical system under study. The optic system consists of a quasi-distributed fiber optic sensor which is based on wavelength-division multiplexing (WDM) and twin-grating sensors. The sensing system has five fundamental components: an optical broadband source, an optical circulator 50/50, an optical spectrometer analyzer (OSA spectrometer), a personal computer, and a quasi-distributed sensor. In particular, the quasi-distributed sensor

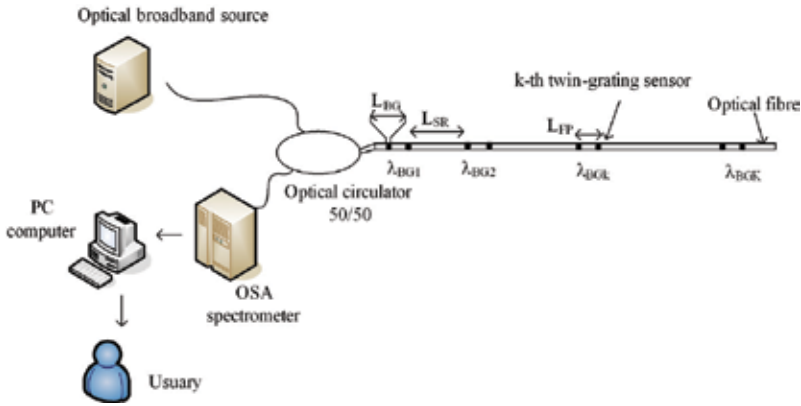


Figure 2. A quasi-distributed sensor based on the WDM technique and twin-grating sensors.

consists of a serial array of twin-grating sensors. Each twin-grating sensor acts as a low-finesse Fabry-Perot interferometer [3, 10]. All interferometers have the same cavity length L_{FP} but each twin-grating sensor has its own Bragg wavelength. Thus, the WDM technique is generated and the FDM technique is eliminated. All Bragg gratings have the same length L_{BG} and a typical reflectivity of 1%. The low reflectivity eliminates cross-talk noise.

2.1. Optical signal

In the sensing system presented in **Figure 2**, the signal from each local sensor is returned by reflection from each twin-grating sensor, where each twin-grating interferometer has its own wavelength. The signal returned to the detector is monitored with the OSA spectrometer; the intensity at each wavelength corresponds to the measurement each local sensor. When the quasi-distributed sensor does not have external perturbations and interference patterns have small variation, the optical signal will be

$$R_T(\lambda) = \sum_{k=1}^K 2a_k \left[\left(\frac{\pi n_1 L_{BG}}{\lambda_{BGk}} \right)^2 \text{sinc}^2 \left(\frac{2n_1 L_{BG} (\lambda - \lambda_{BGk})}{\lambda_{BGk}^2} \right) \right] \left[1 + \cos \left(\frac{4\pi n L_{FP} (\lambda - \lambda_{BGk})}{\lambda_{BGk}^2} \right) \right] \quad (1)$$

The signal parameters are: $R_T(\lambda)$ is a set of interference patterns, λ is the wavelength, λ_{BGk} is the k th Bragg wavelength, a_k are amplitude factors, n_1 is the amplitude of the effective refractive index modulation of the gratings, L_{BG} is the length of gratings, n is the effective index of the core, L_{FP} is the cavity length, and k is the number of twin-grating sensors. The optical signal has the next characteristics: the enveloped function is a sinc one and its width Δ_{BGk} is defined as the spectral distance between its +1 and -1 zeros,

$$\Delta_{BG1} \neq \Delta_{BG2} \neq \Delta_{BG3} \neq \dots \neq \Delta_{BGk} = \frac{\lambda_{BGk}^2}{n_1 L_{BG}} \neq \dots \neq \Delta_{BGK} \quad (2)$$

Each interference pattern has its own central Bragg wavelength and the next condition is true

$$\lambda_{BG1} \neq \lambda_{BG2} \neq \dots \neq \lambda_{BGk} \neq \dots \neq \lambda_{BGK} \quad (3)$$

Interference patterns have approximately the same frequency,

$$\nu_{FP1} \approx \nu_{FP2} \approx \dots \approx \nu_{FPk} = \frac{2nL_{FP}}{\lambda_{BGk}^2} \approx \dots \approx \nu_{FPK} \quad (4)$$

From Eqs. (1) and (4), the cavity length defines the frequency of all interference patterns. Its size can be found in the interval of [3]

$$L_{FPmin} \leq L_{FP} \leq L_{FPmax} \rightarrow 2L_{BG} \leq L_{FP} \leq \frac{\lambda_{BG1}^2}{4n\Delta\lambda} \quad (5)$$

where $\Delta\lambda$ is the spectrometer resolution, $L_{FPmin} = 2L_{BG}$ is the minimum cavity length, and $L_{FPmax} = \frac{\lambda_{BG1}^2}{4n\Delta\lambda}$ is the maximum cavity length. The minimum cavity length is delimited because the Fourier domain phase analysis (FDPA) algorithm does not accept additional information or loss of information. The maximum cavity length is delimited because the OSA spectrometer has a limit of full-width half-maximum (FWHM). **Figure 3** illustrates the optical spectrum.

To know the frequency spectrum $R(\nu)$, the Fourier transform is applied to Eq. (1),

$$R_T(\nu) = \int_{-\infty}^{\infty} R_T(\lambda) e^{-i2\pi\lambda\nu} d\lambda \quad (6)$$

Substituting Eq. (1) into Eq. (6), the spectra $R_T(\nu)$ is

$$R(\nu) = \int_{-\infty}^{\infty} \sum_{k=1}^K 2a_k \left[\left(\frac{\pi n_1 L_{BG}}{\lambda_{BGk}} \right)^2 \text{sinc}^2 \left(\frac{2n_1 L_{BG} (\lambda - \lambda_{BGk})}{\lambda_{BGk}^2} \right) \right] \left[1 + \cos \left(\frac{4\pi n L_{FP} (\lambda - \lambda_{BGk})}{\lambda_{BGk}^2} \right) \right] e^{-i\omega t} dt \quad (7)$$

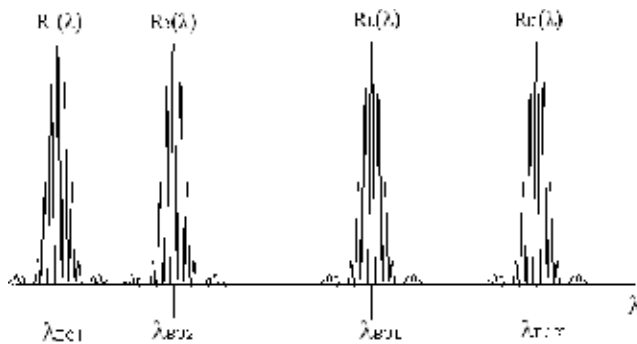


Figure 3. Optical signal detected by the optical analyzer spectrometer.

Solving the transformation, the frequency spectrum is defined by

$$R_T(\nu) = \sum_{k=-K}^K R_k(\nu) = \sum_{k=-K}^K c_k \text{tri}\left(\frac{\nu - \nu_{FPk}}{\nu_{BGk}}\right) \quad (8)$$

This frequency spectrum is the superposition of a set of triangle functions, where a triangle function is defined as $\text{tri}(x) = \begin{cases} 1 - |x| & |x| \leq 1 \\ 0 & \text{otherwise} \end{cases}$, $R_k(\nu)$ is the Fourier transform of the k -th interference pattern, c_k are amplitude factors, ν_{FPk} is the center position of each peak and ν_{BGk} is the bandwidth of each peak.

$$\nu_{BGk} = \frac{4n_1 L_{BG}}{\lambda_{BGk}^2} \quad (9)$$

In the frequency spectrum, the component ν_{FP0} contains information from all twin-grating sensors, the positive components $\nu_{FP1}, \dots, \nu_{FPK}$; the negative components $-\nu_{FP1}, \dots, -\nu_{FPK}$ contain the same information. The minimum bandwidth ν_{BGmin} is

$$\nu_{BGmin} = \frac{4n_1 L_{BG}}{\lambda_{BGK}^2} \quad (10)$$

and the maximum bandwidth ν_{BGmax} is

$$\nu_{BGmax} = \frac{4n_1 L_{BG}}{\lambda_{BG1}^2} \quad (11)$$

Figure 4 shows the frequency spectrum $R(\nu)$. Based on **Figure 4**, all twin-grating interferometers produce approximately the same frequency components. This is possible because all interferometers have the same cavity length and all enveloped functions are approximately similar.

2.2. Optical signal produced by external perturbation

When the quasi-distributed sensor has external perturbations due to the temperature or strain, Bragg gratings and cavity length have an elongation. In turn, interference patterns have a small shift in response to a measured variation. The optical signal detected by the OSA spectrometer is

$$R_T(\lambda, \delta\lambda) = \sum_{k=1}^K 2a_k \left[\left(\frac{\pi n_1 L_{BG}}{\lambda_{BGk}} \right)^2 \text{sinc}^2 \left(\frac{2n_1 L_{BG} (\lambda - \lambda_{BGk} - \delta\lambda_k)}{\lambda_{BGk}^2} \right) \right] \left[1 + \cos \left(\frac{4\pi n L_{FP} (\lambda - \lambda_{BGk} - \delta\lambda_k)}{\lambda_{BGk}^2} \right) \right] \quad (12)$$

It can also be expressed as

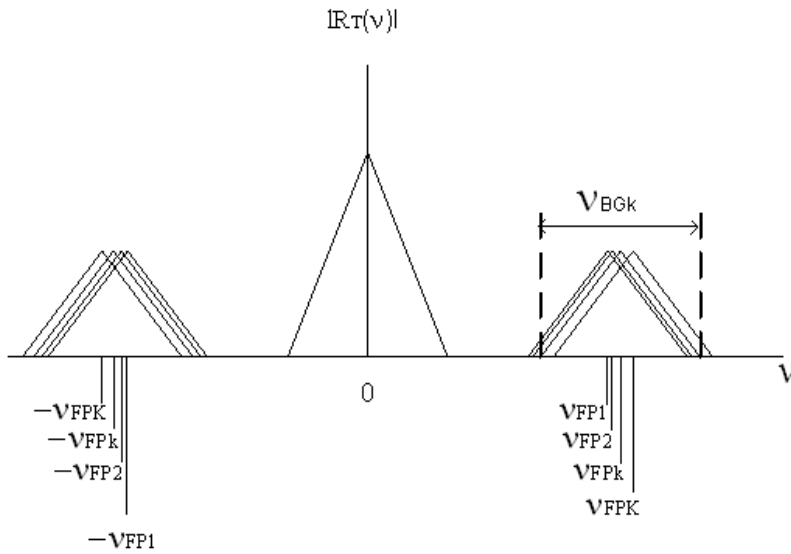


Figure 4. Frequency spectrum determined from the optical signal.

$$R_T(\lambda, \delta\lambda) = \sum_{k=1}^K R_k(\lambda - \delta\lambda_k) = R_1(\lambda - \delta\lambda_1) + \dots + R_k(\lambda - \delta\lambda_k) + \dots + R_K(\lambda - \delta\lambda_K) \quad (13)$$

where $R_T(\lambda, \delta\lambda)$ is the optical signal due to external perturbations and $\delta\lambda_k$ is the Bragg wavelength shift due to measured change [11, 12]. Its frequency spectrum $R_T(v, \delta\lambda)$ is

$$R_T(v, \delta\lambda) = \int_{-\infty}^{\infty} R_T(\lambda, \delta\lambda) e^{-i2\pi\lambda v} d\lambda \quad (14)$$

Substituting Eq. (13) into Eq. (14), the spectra is now

$$R_T(v, \delta\lambda) = \int_{-\infty}^{\infty} \sum_{k=-K}^K R_k(\lambda - \delta\lambda_k) e^{-i2\pi\lambda v} d\lambda \quad (15)$$

Using the Fourier transform properties and solving, the frequency spectra $R_T(v, \delta\lambda)$ takes the form

$$R_T(v, \delta\lambda) = \sum_{k=-K}^K e^{-i2\pi v \delta\lambda_k} R_k(v) \quad (16)$$

$R_T(v, \delta\lambda)$ is the multiplication between $R(v)$ and a set of phases. Each phase contains information of each twin-grating sensor and then the FDPA algorithm can be applied in the demodulation signal.

3. Number of samples

In Ref. [11] the twin-grating sensor was applied for the temperature measurement. The wavelength shift sensitivity to a temperature change was estimated to be $0.00985 \text{ nm}/^\circ\text{C}$. The demodulation signal was done using the Fourier domain phase analysis algorithm. The optical signal was acquired applying direct spectrometric detection. This detection technique uses an optical spectrometer analyzer; then, the acquired optical signal becomes discrete. The signal samples $R_T(\lambda_j)$ are taken as wavelengths $\lambda_j = \lambda_{\min} + j\delta\lambda_s$, where $j = 0, 1, \dots, N - 1$, N is the number of samples. The interval working is $\lambda_w = \lambda_{\max} - \lambda_{\min}$: λ_{\max} is the maximum wavelength, λ_{\min} is the minimum wavelength and $\delta\lambda_s$ is the wavelength step.

From **Figure 4**, the maximum frequency ν_{\max} is

$$\nu_{\max} = \nu_{FPK} + \frac{1}{2}\nu_{BG\max} \quad (17)$$

Substituting Eqs. (4)–(11) into Eq. (17), the maximum frequency is

$$\nu_{\max} = \frac{2nL_{FP\max}}{\lambda_{BG1}^2} + \frac{2n_1L_{BG}}{\lambda_{BG1}^2} \quad (18)$$

When we substitute the maximum cavity length (Eq. (5)) into Eq. (18), the parameter ν_{\max} takes the form

$$\nu_{\max} = \frac{1}{2\Delta\lambda} + \frac{2n_1L_{BG}}{\lambda_{BG1}^2} \quad (19)$$

Applying the sampling theorem, the sampling frequency ν_s is

$$\nu_s \geq 2\nu_{\max} = \frac{1}{\Delta\lambda} + \frac{4n_1L_{BG}}{\lambda_{BG1}^2} \quad (20)$$

Since $\nu_s = \frac{1}{\delta\lambda_s}$, we have

$$\delta\lambda_s \leq \frac{\Delta\lambda\lambda_{BG1}^2}{\lambda_{BG1}^2 + 4n_1L_{BG}\Delta\lambda} \quad (21)$$

Finally, the number of samples N is given by

$$N = \frac{\lambda_w}{\delta\lambda_s} = \frac{\lambda_w(\lambda_{BG1}^2 + 4n_1L_{BG}\Delta\lambda)}{\Delta\lambda\lambda_{BG1}^2} \quad (22)$$

Samples N depend on twin-grating sensor properties, the spectrometer resolution, and the interval working. The number of samples is a very important parameter for the twin-grating sensor demodulation because it affects the sensor's resolution.

4. Capacity of wavelength-division multiplexing

In Refs. [4, 5] two experimental sensing systems where twin-grating fiber optic sensors were applied on wavelength-division multiplexing were reported. The first optical system consisted of two wavelength channels. Both channels were centered around 815 and 839 nm. The second optical system consisted of three wavelength channels. The channels were around 1542, 1548, and 1554 nm. Therefore, based on the Bragg grating characteristics, the twin-grating interferometer can be applied in wavelength-division multiplexing if and only if each interferometer sensor has its own Bragg wavelength: $\lambda_{BGk} = 2n_1\Lambda_k$, where Λ_k is the period [1]. In this case, each interference pattern has its own bandwidth in the wavelength domain. The patterns are in the interval of λ_{min} until λ_{max} (interval working λ_w); it is not possible in other positions, see **Figure 5**.

Let us introduce the operation range $\Delta\lambda_{op}$; the operation range defines the interval in which an interference pattern can move into the wavelength domain. Each interferometer sensor has its own operation range and overlapping is not acceptable. To calculate the capacity of wavelength-division multiplexing K , we use the interval working λ_w and the operation range $\Delta\lambda_{op}$

$$K = \frac{\lambda_w}{\Delta\lambda_{op}} = \frac{\lambda_{max} - \lambda_{min}}{\Delta\lambda_{op}} \tag{23}$$

K is the number of local sensors and wavelength channels. λ_w , λ_{min} , λ_{max} , and $\Delta\lambda_{op}$ parameters can be observed in **Figure 3**.

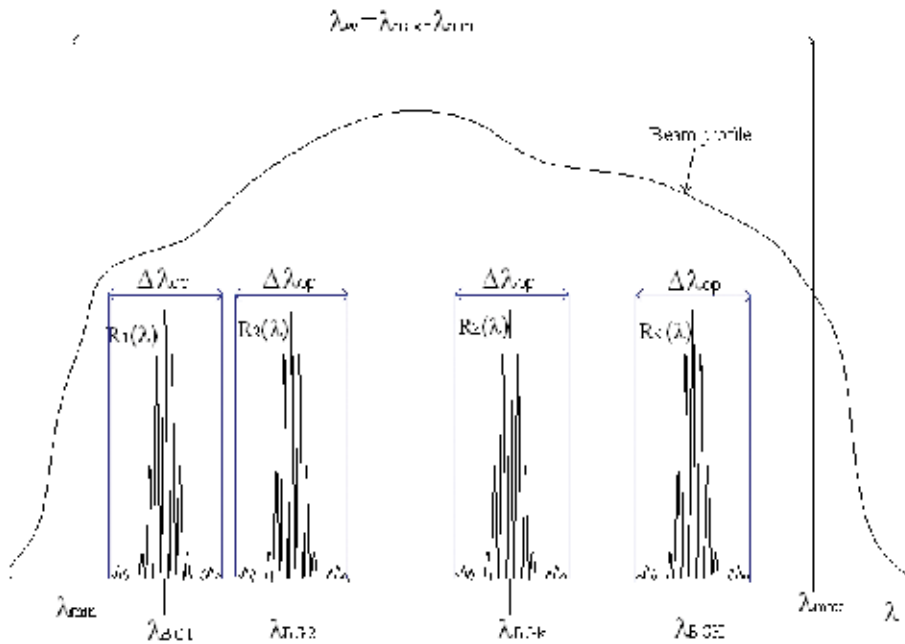


Figure 5. λ_w , λ_{min} , λ_{max} , and $\Delta\lambda_{op}$ representation.

To illustrate, the next numerical example is presented.

Example 1: A broadband light source has the interval from $\lambda_{\min} = 1470$ nm to $\lambda_{\max} = 1620$ nm. If the operation range is selected to $\Delta\lambda_{\text{op}} = 6$ nm, the number of local sensors is $K = \frac{1620-1470}{6} = 25$. The quasi-distributed sensor would have 25 twin-grating sensors and the signal will have 25 wavelength channels. Here, two important points can be mentioned: 1) The parameter $\Delta\lambda_{\text{op}}$ permits the selection of the number of local sensors; 2) the cost per sensing point can be reduced combining the wavelength-and-frequency division multiplexing.

5. Demodulation signal

In this section, we present the demodulation signal for the quasi-distributed sensor based on wavelength-division multiplexing. The signal processing combines the Fourier domain phase analysis (FDPA) algorithm, a bank of K filters and a band-pass filter. The FDPA algorithm was described and also applied in Refs. [3, 11]. The bank of K filters can be defined as

$$F(\lambda) = \text{rect}\left(\frac{\lambda}{\Delta\lambda_{\text{op}}}\right) \otimes \sum_{k=1}^K \delta(\lambda - \lambda_{BGk}) \quad (24)$$

where the symbol \otimes indicates the convolution operation, the rect function (Eq. (24)) is defined as

$$\text{rect}(\lambda) = \begin{cases} 1 & |\lambda| < \frac{\Delta\lambda_{\text{op}}}{2} \\ 0 & |\lambda| > \frac{\Delta\lambda_{\text{op}}}{2} \end{cases} \quad (25)$$

and δ is the Dirac delta. Invoking the Dirac delta properties, the bank of K filters is

$$F(\lambda) = \sum_{k=1}^K \text{rect}\left(\frac{\lambda - \lambda_{BGk}}{\Delta\lambda_{\text{op}}}\right) \quad (26)$$

The signal $F(\lambda)$ is a series of rect functions in the wavelength domain; its bandwidth is the operation range and the central positions are the Bragg wavelengths. On the other hand, $R_k(\nu)$ is the Fourier transform for the k th interference pattern. The spectrum $R_k(\nu)$ consists of three triangle functions. The component ν_{FP0} contains information from all twin-grating sensors and this signal cannot be used in the demodulation signal. The ν_{FPk} and $-\nu_{FPk}$ components contain the same information and any component can be used in the demodulation signal. Then, the band-pass filter can be defined as

$$F(\nu) = \text{rect}\left(\frac{\nu - \nu_{FPk}}{\nu_{BGk}}\right) \quad (27)$$

where the rect function (Eq. (27)) has the next definition

$$rect(v) = \begin{cases} 1 & |v| < \frac{v_{BGk}}{2} \\ 0 & |v| > \frac{v_{BGk}}{2} \end{cases} \quad (28)$$

The filter $F(v)$ is a rect function in the frequency domain: its bandwidth is v_{BGk} and its central position is v_{FPk} .

Basically, the digital demodulation consists of two stages: calibration and measurement. The calibration stage is developed once; five steps are necessary and the references are generated. The calibration considers the signal acquisition $R_T(\lambda)$, filtering signal $R_k(\lambda)$, separation of frequency components $R_k(v)$, filtering $R_m(\lambda)$, and its complex conjugate $R_m^*(\lambda)$, where the symbol * indicates the complex conjugate. The measurement stage is developed for each measurement and eight steps are necessary. The measurement stage considers the signal acquisition $R_T(\lambda, \delta\lambda)$, filtering $R_k(\lambda, \delta\lambda)$, separation of frequency components $R_T(v, \delta\lambda)$, filtering $R_m(v, \delta\lambda)$, comparison between spectrums $R_m^*(\lambda)$ and $R_m(v, \delta\lambda)$, $2\pi P$ ambiguity elimination; Bragg wavelength shift evaluation and an adaptive filter is applied. The adaptive filter is a set of coefficients as was described in Ref. [13]. The complete procedure is shown at **Figure 6**.

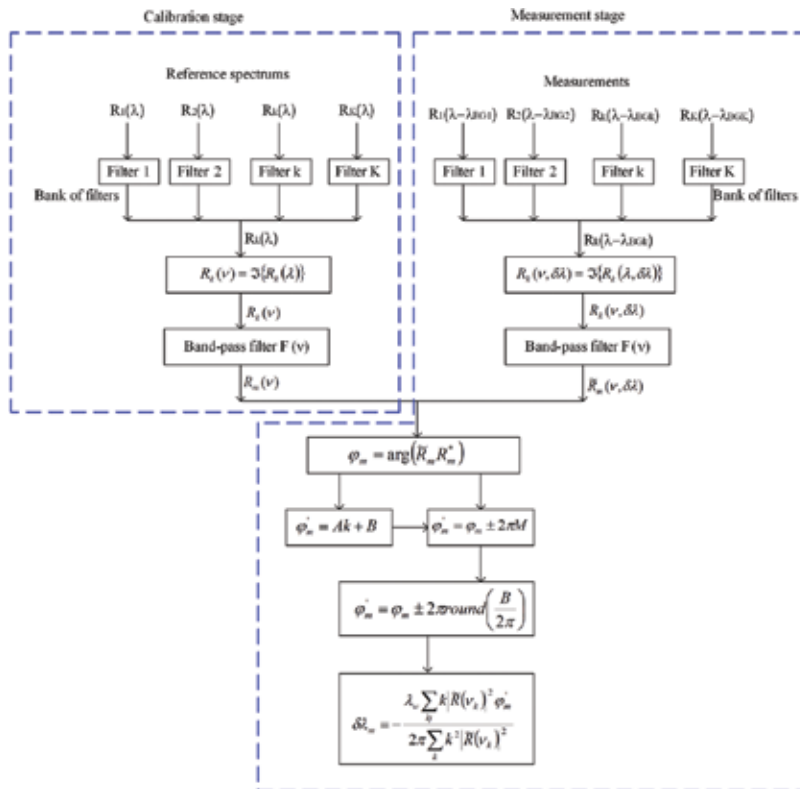


Figure 6. Digital demodulation represented schematically: $R_m^* = R_m^*(\lambda)$, $\bar{R}_m = \bar{R}_m(v, \delta\lambda)$ and the symbol * indicates the complex conjugate.

6. Numerical simulation and discussion

6.1. Parameters and results

In Ref. [3], the twin-grating fiber optic sensor was applied on frequency-division multiplexing. The numerical results confirmed that the twin-grating sensor has a high resolution and the resolution is a function of cavity length. In the numerical simulation, the number of samples was $N = 1024$, the noise was in the interval $\sqrt{SNR} = 10^0$ to $\sqrt{SNR} = 10^4$. The quasi-distributed sensor parameters and optical signal parameters can be observed in **Table 1**.

Being aware of that, our goal is to apply the twin-grating interferometer to wavelength-division multiplexing, a quasi-distributed fiber optic sensor is numerically simulated as was done in Ref. [3]. The quasi-distributed sensor consists of three twin-grating sensors. The physical parameters are shown in **Table 2**. In our numerical simulation, we use some parameters from Ref. [3]: $L_{FP1} = L_{FP2} = L_{FP3} = L_{FP4} = 4$ (mm), $L_{BG} = 0.5$ (mm), $n = 1.46$, $N = 1024$; the Bragg gratings have rectangular profile; the noise has Gaussian distribution and its value is in the interval $\sqrt{SNR} = 10^0$ to $\sqrt{SNR} = 10^4$. The same values are possible because our goal is to prove the wavelength-division multiplexing. The measurements are in the intervals $S1 \rightarrow -0.3$ to 0.3 nm, $S2 \rightarrow -0.5$ to 0.5 nm and $S3 \rightarrow -1$ to 1 nm. The signal parameters are shown in **Table 2** while **Figure 5** shows our numerical results: Demodulation errors vs. $SNR^{1/2}$. A Laptop Toshiba 45C was used with a 512 Mb RAM and a velocity of 1.7 GHz.

Analyzing **Tables 1** and **2**, the sensing system presented in **Figure 3** is based on wavelength-division multiplexing where their wavelength-channels are 1531.5, 1535.5, and 1539.5 nm. The frequency-division multiplexing was eliminated since the frequency components are $v_{FP1} = 4.979$ (cycles/nm), $v_{FP2} = 4.953$ (cycles/nm), and $v_{FP3} = 4.928$ (cycles/nm). All interference patterns have approximately the same bandwidth: $\Delta\lambda_{BG1} = 3.213$ (nm), $\Delta\lambda_{BG2} = 3.229$ (nm), and

Sensor number	Sensor parameters	Signal values
Twin-grating sensor 1 (S1)	$L_{FP1} = 4$ (mm)	$\Delta\lambda_{BG} = 3.22$ (nm)
	$L_{BG} = 0.5$ (mm)	$v_{FP1} = 4.95$ (ciclos/nm)
	$n = 1.46$	$v_{BG} = 1.23$ (ciclos/nm)
	$\lambda_{BG} = 1532.5$ (nm)	
Twin-grating sensor 2 (S2)	$L_{FP2} = 8$ (mm)	$\Delta\lambda_{BG} = 3.22$ (nm)
	$L_{BG} = 0.5$ (mm)	$v_{FP2} = 9.91$ (ciclos/nm)
	$n = 1.46$	$v_{BG} = 1.23$ (ciclos/nm)
	$\lambda_{BG} = 1532.5$ (nm)	
Twin-grating sensor 3 (S3)	$L_{FP3} = 16$ (mm)	$\Delta\lambda_{BG} = 3.22$ (nm)
	$L_{BG} = 0.5$ (mm)	$v_{FP3} = 19.82$ (ciclos/nm)
	$n = 1.46$	$v_{BG} = 1.23$ (ciclos/nm)
	$\lambda_{BG} = 1532.5$ (nm)	

Table 1. Sensor parameters and signal parameters used in the frequency-division multiplexing [3].

Sensor number	Sensor parameters	Signal values
Fabry-Perot sensor 1 (S1)	$L_{FP1} = 4$ (mm)	$\nu_{FP1} = 4.979$ (cycles/nm) (Eq. (2))
	$L_{BG} = 0.5$ (mm)	$\Delta_{ABG1} = 3.213$ (nm) (Eq. (3))
	$n = 1.46$	$\nu_{BG1} = 1.244$ (cycles/nm) (Eq. (9))
	$\lambda_{BG1} = 1531.5$ (nm)	
Fabry-Perot sensor 2 (S2)	$L_{FP2} = 4$ (mm)	$\nu_{FP2} = 4.953$ (cycles/nm) (Eq. (2))
	$L_{BG} = 0.5$ (mm)	$\Delta_{ABG2} = 3.229$ (nm) (Eq. (3))
	$n = 1.46$	$\nu_{BG2} = 1.238$ (cycles/nm) (Eq. (9))
	$\lambda_{BG2} = 1535.5$ (nm)	
Fabry-Perot sensor 3 (S3)	$L_{FP3} = 4$ (mm)	$\nu_{FP3} = 4.928$ (cycles/nm) (Eq. (2))
	$L_{BG} = 0.5$ (mm)	$\Delta_{ABG3} = 3.246$ (nm) (Eq. (3))
	$n = 1.46$	$\nu_{BG3} = 1.232$ (cycles/nm) (Eq. (9))
	$\lambda_{BG3} = 1539.5$ (nm)	

Table 2. Quasi-distributed sensor parameters used and signal values obtained.

$\Delta\lambda_{BG3} = 3.246$ (nm). Finally, the triangle functions have approximately the same bandwidth: $\nu_{BG1} = 1.244$ (cycles/nm), $\nu_{BG2} = 1.238$ (cycles/nm) and $\nu_{BG3} = 1.232$ (cycles/nm). Thus, we prove that the quasi-distributed sensor applies the wavelength-division multiplexing and the FDM technique was eliminated.

Figure 7 Shows the behavior of demodulation error vs. signal-to-noise rate $SNR^{1/2}$, if demodulation error is denominated as the resolution. These results confirm good resolution for the twin-grating sensors. Since the FDPA algorithm makes two evaluations of the Bragg wavelength shift [3, 7], the twin-grating sensors have two resolutions: Low resolution σ_{env} and high resolution σ_{int} . the threshold [3].

$$\sigma_{env} < \frac{\lambda_{BGk}^2}{12nL_{FP}} \tag{29}$$

specifies the boundary between low resolution and high resolution. Using Eq. (29) and **Table 2**, the thresholds are approximately $S1 \rightarrow \sigma_{env} < 0.0335$ nm, $S2 \rightarrow \sigma_{env} < 0.0336$ nm, and $S3 \rightarrow \sigma_{env} < 0.0338$ nm. For each local twin-grating sensor, an $SNR^{1/2}$ threshold is observable:

$S1(1531.5 \text{ nm}) \rightarrow SNR^{\frac{1}{2}} \approx 10^{1.15}$, $S2(1535.5 \text{ nm}) \rightarrow SNR^{\frac{1}{2}} \approx 10^{1.19}$, $S3(1539.5 \text{ nm}) \rightarrow SNR^{1/2} \approx 10^{1.2}$. The threshold values are very close because the twin-grating sensors have similar cavity length and the enveloped functions have approximately the same bandwidth $\Delta\lambda_{BGk}$.

6.2. Discussion

The quasi-distributed fiber optic sensor (**Figure 3**) would be built on wavelength-division multiplexing and twin-grating interferometers. Our results optimize the sensor’s implementation and also permit its design. Local sensor properties, light source characteristics, noise (Gaussian distribution), signal processing and detection technique are considered in our

numerical simulation. Our experimental results (**Figure 7**) corroborate well functionally. Two resolutions are also confirmed σ_{BG} for each local sensor: low resolution σ_{env} and high resolution σ_{int} . Both resolutions depend on noise system, cavity length, instrumentation, sensor properties, and the digital demodulation algorithm.

A twin-grating fiber optic sensor and an optical fiber sensor based on a single Bragg grating will have the same resolution if and only if FDPA algorithm cannot eliminate the $2\pi P$ ambiguity. In this case, the twin-grating sensor has only low resolution because the FDPA algorithm evaluates the Bragg wavelength shift with the enveloped function [11]. However, if the optical sensing system has small noise, then the twin-grating sensor has high resolution, since the Bragg wavelength shift was evaluated combining the enveloped and modulated functions [3, 11]. Additionally, three twin-grating sensors have approximately the same resolution because three wavelength channels are very close and all interferometer sensors have the same cavity length. In conclusion, the sensor's resolution is high while the FDPA algorithm can acceptably demodulate the optical signal.

The presented study optimizes the quasi-distributed sensor which was shown in **Figure 3**. Combining our study (this work) and the analysis presented in Ref. [3], the experimental sensing system described by Shlyagin et al. [4, 5] can be optimized. The optimization will be on signal processing, local sensor properties, sensitivity, resolution and instrumentation parameters. Additionally, the cost per sensing point is considerably reduced.

Our future work has the following direction: a theoretical analysis and practical application. In the theoretical analysis, frequency-and-wavelength division multiplexing can be implemented based on the twin-grating interferometer; resolution is another direction. In the practical

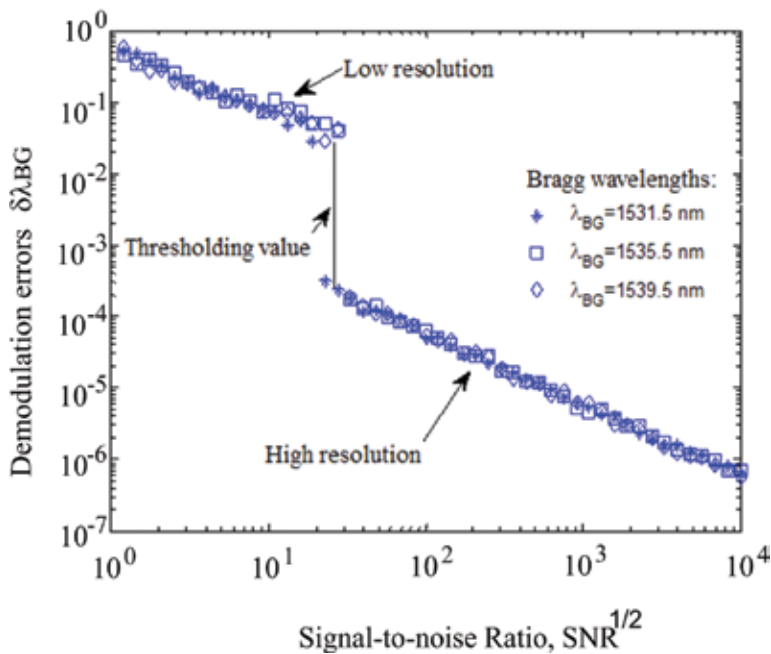


Figure 7. Numerical results obtained from the numerical experiments.

applications, the quasi-distributed sensor can be applied for temperature monitoring, gasoline detection (security), strain measurement, and level liquid measurement. Our analysis makes an excellent contribution to quasi-distributed sensor implementation because all local sensors will have high resolution (see **Figure 7**), high sensibility, low cost by sensing point, and the quasi-distributed sensor can be designed without other requirements.

7. Conclusion

In this work, a quasi-distributed fiber optic sensor was numerically simulated. The sensor was based on twin-grating sensors and wavelength-division multiplexing. The numerical results show the resolution for each local twin-grating sensor. Local sensors have approximately the same resolution because all twin-grating sensors have the same cavity length and the wavelength channels are close. Two resolutions were obtained for each local sensor. Our numerical results show that the quasi-distributed sensor has potential industrial application: temperature measurement, strain measurement, pressure measurement, humidity monitoring, and security system.

Acknowledgements

Authors thank PRODEP 2017 No. F-PROMEPE-39/Rev-04 SEP-23-005 (number DSA/103.5/16/10313), PRODEP 2017 Project No. 236110 of found 1.1.9.25 (Agreement RG/003/2017) and PRODEP 2017 Project No 238635 (511-6/17-8091).

Author details

José Trinidad Guillen Bonilla^{1,2*}, Héctor Guillen Bonilla³, Antonio Casillas Zamora¹, Gustavo Adolfo Vega Gómez¹, Nancy Elizabeth Franco Rodríguez⁴, Alex Guillen Bonilla⁵ and Juan Reyes Gómez⁶

*Address all correspondence to: trinidad.guillen@academicos.udg.mx

1 Electronic Department, CUCEI, University of Guadalajara, Guadalajara, Jalisco, México

2 Mathematics Department, CUCEI, University of Guadalajara, Guadalajara, Jalisco, México

3 Department of Engineering Projects, CUCEI, University of Guadalajara, Guadalajara, Jalisco, México

4 Computer Department, CUCEI, University of Guadalajara, Guadalajara, Jalisco, México

5 Department of Computer Science and Engineering, CUValles, University of Guadalajara, Ameca, Jalisco, México

6 Faculty of Chemical Science, Colima University, Coquimatlán, Colima, Mexico

References

- [1] Kashyap R. Photosensitive optical fibers: Device and applications. *Optical Fibre Technology*. 1994;1:17-34
- [2] Grattan KTV, Sun DT. Fiber optic sensor technology: An overview. *Sensors and Actuators*. 2000;82:40-61
- [3] Guillen Bonilla JT, Guillen Bonilla A, Rodríguez Betancourt VM, Guillen Bonilla H, Casillas Zamora A. A theoretical study and numerical simulation of a quasi-distributed sensor based on low-finesse Fabry-Perot interferometer: Frequency-division multiplexing. *Sensors*. 2017;17:859. DOI: 10.3390/s17040859
- [4] Shlyagin MG, Miridonov SV, Márquez-Borbón I, Spirin VV, Swart PL, Chitchebakov AA. Multiplexed twin Bragg grating interferometer sensor. In: *Proceedings of the Optical Fiber Sensors Conference Technical Digest (OFS 2002)*; 10 May 2002; Portland, OR, USA. pp. 191-194
- [5] Shlyagin MG, Miridonov SV, Tentori Santa-Cruz D, Mendieta Jiménez FJ, Spirin VV. Multiplexing of grating-based fiber sensors using broadband spectral coding. In: *Proceedings of the Conference on Fiber Optic and Laser Sensors and Applications*; 1 November 1998; Boston, MA, USA; Volume 3541, pp. 271-277
- [6] Cibula E, Donlagic D. In-line short cavity Fabry-Perot strain sensor for quasi-distributed measurement utilizing standard OTDR. *Optics Express*. 2007;15(14):8728
- [7] Yu XJ, Zhang YL, Li K, Zhang JT, Lv HJ, Liu SC. Frequency-division multiplexing sensing system based on multilongitudinal mode fiber lasers and beat frequency demodulation. *IEEE Photonics Journal*. 2015;7(2):7. ID: 7901307. DOI: 10.1109/JPHOT.2015.2418265
- [8] Wang Y, Gong J, Wang DY, Dong B, Bi W, Wang A. A quasi-distributed sensing network with time-division-multiplexing fiber Bragg gratings. *IEEE Photonics Technology Letters*. 2011;23(2):70-72. DOI: 10.1109/LPT.2010.2089676
- [9] Yuan L, Zhou L, Jin W, Yang J. Design of a fiber-optic quasi-distributed strain sensors ring network based on a white-light interferometric multiplexing technique. *Applied Optics*. 2002;41(34):7205-7211. DOI: 10.1364/AO.41.007205
- [10] Martinez-Manuel R, Shlyagin MG, Miridonov SV, Meyer J. Vibration disturbance localization using a serial array of identical low-finesse fiber Fabry-Perot interferometers. *IEEE Sensors*. 2012;12(1):124-127. DOI: 10.1109/JSEN.2011.2119479
- [11] Miridonov SV, Shlyagin MG, Tentori D. Twin-grating fiber optic sensor demodulation. *Optics Communication*. 2001;191:253-262

- [12] Miridonov SV, Shlyagin MG, Tentori D. Digital demodulation of twin-grating fiber-optic sensor. In: Proceedings of the Conference on Fiber Optics and Laser Sensors and Applications, Boston, MA, USA; 1 November 1998; Volume 3541; pp. 33-40
- [13] Miridonov SV, Shlyagin MG, Spirin VV. Resolution limits and efficient signal processing for fiber optic Bragg grating with direct spectroscopic detection. In: Proceedings of the Conference on Optical Measurement Systems for Industrial Inspection III, Munich, Germany; 2003; Volume 5144; pp. 679-686

Numerical Modeling of Chemical Compounds' Fate and Kinetics in Living Organisms: An Inverse Numerical Method for Rate Estimation from Concentration

Kovacs Eموke Dalma and Kovacs Melinda Haydee

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.76611>

Abstract

Emerging chemical compounds are ubiquitous in all environmental compartments and may pose a risk to biota ecosystems. The quantification and prediction of environmental partitioning of these chemicals in various environmental compartment systems (water, sediments, soil, air, biota) is an important step in the comprehensive assessment of their sources, fates, and not finally of their uptake potential by various living organisms of ecosystems.

Any numerical solution that has as a final goal "prediction" requires a large number of experimental data. In case of environmental studies of chemical compounds, monitoring most studies is costly, time-consuming, and requires both qualified personnel and high-precision equipment. Finding a suitable numerical model that could predict the fate of chemicals could be extremely useful, facilitating those environmental scientists, users, managers, authorities, and corresponding decision-makers for a more conscious use of these substances, thus protecting the environment and biota.

Considering the mentioned disadvantages regarding chemical compounds' monitoring, the aim of this research is to find numerical solutions that enable the prediction of such chemical compounds' fate under different environmental compartments and the uptake potential by living organisms as plants. The concept of the inverse numerical method was used in order to find chemical compounds' rate of accumulation in various environmental matrixes and potential uptake by living organisms, all starting from the chemical compounds' concentrations.

Keywords: numerical modeling, simulation, prediction, pollutants fate, kinetics

1. Introduction

Thousands of chemicals are used in industry, agriculture, pharmacy, commerce, and daily life. With that, a large number of chemical compounds enter the environment. Often these are considered with potential harmful effects on environmental media quality and biota safety. Thus, the understanding of them is crucial both for a better management of their use and for the better protection of the environment and living organisms. Monitoring these chemical compounds frequently is time-consuming and requires large financial efforts; one of the most cost-effective as well as time-efficient methods of evaluating their behavior in the environment and living organisms could be the use of predictive numerical models.

Experimental data and models for chemical compounds' fate and kinetics in living organisms play a crucial role for assessing the potential human and ecological risks associated with chemical use.

Plants are receptor organisms and could be either direct or indirect vectors for chemical exposure to all other organisms [1]. In the first instance, the generated experimental data considering chemical concentrations in different media of the environment and biota are necessary to improve our understanding on plant-chemical-environment interactions. These, in turn, admit and bring forward the development of better scientific knowledge as well as conceptual and predictive models on chemical partition, fate, and uptake [2]. The strong interconnections between experimental data and model development are continuous and a long-term updated process which is needed to advance our ability to provide reliable quality information that can be used in various environmental protection contexts and regulatory risk assessments [1].

At this moment there are no standard protocols both for chemical compounds' bioaccumulation data generation and for data use for prediction through numerical methods [1, 3]. For the reliable modeling of plant-chemical-environment interactions with the major goal to predict chemical compounds' fates and kinetic in living organisms, it is necessary to understand and keep into account all process, phenomena, and characteristics of both chemicals and receptors (living organisms and the environment) and the interconnected process between them. Inconsistent data collection, inaccurate generation of them, or reporting them with gaps will provide improper and less useful information for their application in assessment and numerical model development.

This chapter is about to find optimal numerical modeling ways considering chemical compounds' fate and kinetics in a living organism, specifically plants. The aim of this study was to propose a numerical procedure which estimates the highest accumulation rate of a chemical compound of interest for a growing living organism and to validate the procedure.

2. Problem formulation

- i. *The presence of chemical compounds in the environment and biota* is the subject of global interest with the general aim to ensure that their impact on humans, other living organisms, and the environment is minimized. Chemicals that once enter the environment can persist or support different types of transformations resulting in new transformation products. Based

on the characteristics of transformation products, these transformation processes could be either of benefit (attenuation processes) or be hazardous (when the resultant transformation products possess more hazardous characteristics than the “parent” compound).

- ii. *Chemical compounds' transformation in the environment:* Based on chemical compounds' specific physicochemical properties as well of biogeochemical and physicochemical characteristics of the media in which they are discharged (soil, water, air), these chemicals can be distributed across different environmental compartments (soil, surface and/or underground water , air, sediments, etc.) and biota [4]. In almost all cases their accumulation in environmental or biota compartments is characterized by continuous dynamic processes such as volatilization, degradation, precipitation, sorption, and so on, processes that often have the potential to end up in the formation of new chemical compounds called as transformation products of the initial chemical compound (**Figure 1**). Often these compounds could enter in new reaction processes and for other new end products. A schematic diagram of principal processes that could take place in both environment and biota and their interconnection with the “initial” chemical compounds and “resultant” transformation products is presented in **Figure 1**.

Challenges considering potential transformation products are given by their physicochemical properties that in several cases could make them more harmful against environment quality or biota safety than the initial chemical compound. Both processes as well as the resultant transformation products' formation are directly dependent on the environmental conditions as well properties of chemicals.

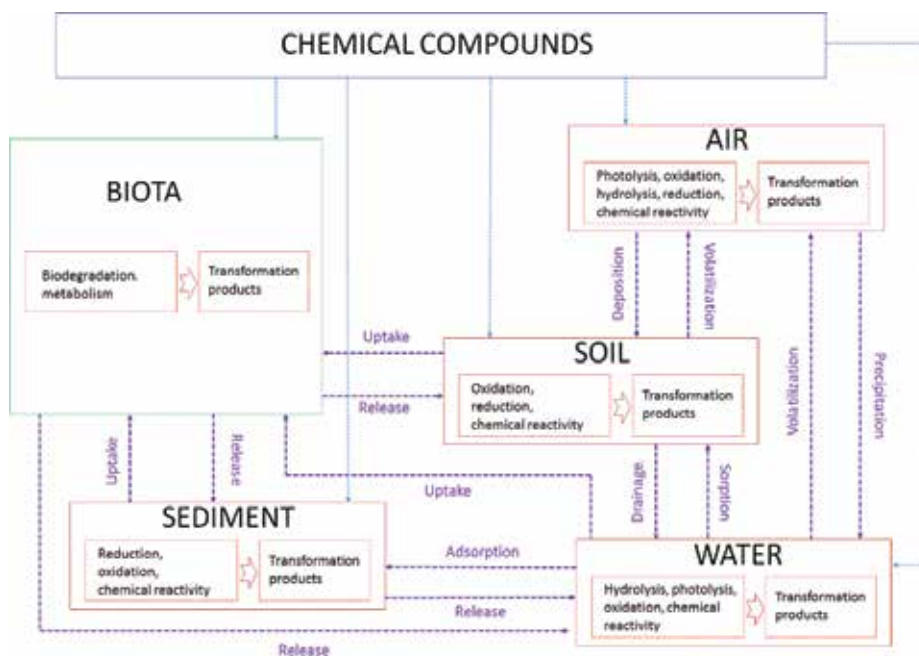


Figure 1. A schematic presentation of the main processes involved in chemical compounds and corresponding potential transformation of products' fate between different environmental compartments and biota.

- i. *Chemical compounds' uptake by plants:* Chemicals' uptake by plant organisms is a system of complex and multi-step processes. These processes could be classified firstly as chemical uptake and transportation between different anatomical compartments (e.g., root to any other anatomical compartment) and secondly as chemical uptake from different environmental compartments (route of exposure) and plant anatomical compartments (particle deposition, vapor uptake from the atmosphere, and so on). The amplitude of these processes is determined by physicochemical properties of the chemical that is under uptake [5]. Current literature presents clear-cut evidence that the availability of most organic chemical compounds is governed on the one hand by their lipophilicity and on the other hand depends on the organic matter (OM) content of the soil under consideration [6]. Some compounds form "bound" residues with organic matter (OM) or humus particles in the soil. Besides, the nature and rooting pattern of the vegetation will have greater influence on the solubility of chemicals. Exuding up to 25% of the net carbon fixed during photosynthesis into the rhizosphere, plants modify given soil-chemical interactions in multiple ways. Secondary plant products (phenolic) and soil bioactive compounds (carbohydrates, organic acids, etc.) could also impact soil micro-biodiversity that could influence in a positive way transformation of organic pollutants to reactive metabolites [7]. For example, it has been demonstrated that isoproturon is metabolized to available plant and reactive compounds in rhizosphere soil [8], while the bacterial conversion of arochlors to reactive metabolites has been one of the early results of bioremediation studies [9].
- ii. Probably, one of the most effective ways to study chemical behavior and fate is to *use mathematical fate models*. Mechanistic environmental models use mathematical equations which describe the parameters of an environment (e.g. data on flows, depths, pH, temperature, etc.) interconnected with the physicochemical properties of the chemical compounds under various conditions with the final aim of predicting their fate in the environment. According to [10], this can be an inexpensive and suitable approach for setting the limits for discharges in the environment of certain chemical compounds, and since the initial parameter description has been set up and validated by in situ and laboratory data, it can be studied with a minimum set of analysis (e.g., only the quantification of chemical compound inputs to the ecosystem) [10].

3. Framework for chemical compound interaction with the environment and living organisms

The ability of numerical models to accurately predict concentrations of target chemical compounds in any living organism depends on the model's ability to mimic the processes involved in their uptake, and this must be assessed before they can be confidently applied [10]. After that it is necessary to consider all of these processes in order to include them in the numerical model that wants to be developed [11].

Soil-root transport: The uptake of chemicals by the root from the soil is mediated in high percentage by soil water content through the plant transpiration process [12]. A large number of organic chemicals also can be sorbed or bound to the components in soil (clay, iron oxides, organic matter),

those often found in the rhizosphere in significant amounts [13, 14]. Also, lipophilic organic chemicals possess a greater tendency to partition into plant roots than hydrophilic chemicals. Although chemical properties are important predictors of the uptake potential, the physiology and composition of the plant root itself is also a significant influence, with differences in the uptake potential explained by the varying types and amounts of lipids in root cells [15]. Uptake from the external media is often expressed as a root concentration factor (RCF), which is the ratio of chemical concentration in the root to the concentration found in external media [13, 16].

Transfer from roots to other anatomical compartments of plants: The major factor that illustrates the amount of a chemical compound that was transferred from the plant root part to other anatomical compartments is the transpiration stream concentration factor (TSCF) which is the ratio of chemical concentration in transpiration to the concentration found in the external part. TSCF could be predicted from knowledge of the chemical compound lipophilicity, with maximum uptake, a $\log K_{OW}$ about 1 [17]. Once the chemical is transported to the stem, plant water and solutes take it and continue to transport these chemicals to the rest of the anatomical compartments through vascular systems and cell tissues [18].

Vapor or gas uptake from ambient air: Another exposure route with chemicals in case of plant materials could be the ambient air that contains a large number of contaminants. This exposure route is governed by gaseous exchange and facilitates the transport and uptake of chemicals that are volatiles and which are more easily partitioned in air than in water. This has been shown to be the main uptake pathway in the above-ground plant parts for a variety of chemical compounds (e.g. PCBs, tetra- and hexa-chlorinated PCDD/Fs) [19]. Previous studies have reported a good correlation between shoot uptake and chemical properties of compounds (K_{OW} Henry's Law constant, octanol-air partition coefficient, etc.). In studies presented by [20], it was evidenced that gaseous uptake is the primary pathway for chemicals with an octanol-air partition coefficient ($\log K_{OA}$) less than 11 [20].

Particulate deposition on plant surfaces: Pesticides as well as other chemical contaminants are bound to soil particles which may be transported by wind and/or rain and deposited on the above-ground anatomical compartments of plants. Studies presented by [5, 21], evidenced that dry deposition onto the leaf of suspended particles that contain PCDD/Fs is the major route of uptake due to PCDD/Fs permeation through the cuticle. Similarly, in their studies, wet deposition was shown which could also be the dominant deposition mechanism for organic chemicals with Henry's Law constant of less than 1×10^{-6} [5, 21].

To date, a number of mathematical models have been developed to facilitate the exposure assessment of chemical contaminants, with important results in the modeling of pollutants' multimedia fate and the modeling of pollutants' linkage with transformation products, especially in water environmental compartments [4].

4. Case study presentation

The properties of wild growing mushrooms make them valuable resources both in culinary practices and in pharmaceutical practices. They are recognized as healthy food with low

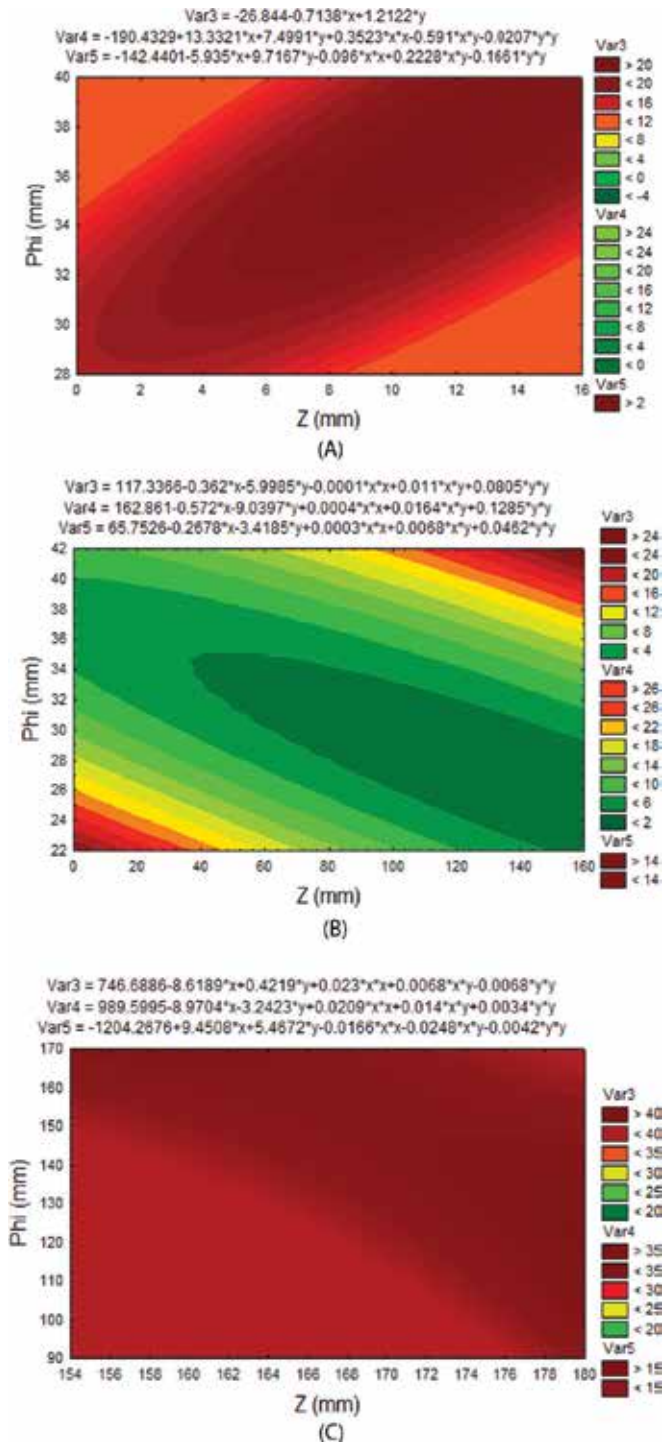


Figure 2. Chemical compound concentration variations in different anatomical compartments: (a) concentration variation in the first anatomical compartment (basal bulb) of the mushroom; (b) concentration variation in the second anatomical compartment (stipe) of the mushroom; (c) concentration variation in the third anatomical compartment (cap) of the mushroom.

contents of calories and fats but high in vitamins, minerals, and vegetable proteins. Their suitability for use by the pharmaceutical industry is given by their rich antioxidant chemical constituents that are capable of preventing the human body from oxidative damage [22]. It is also known that mushrooms could be considered as good bioindicators for the evaluation of environmental pollution, since they are known to accumulate a broad range of chemical compounds [23]. The aim of this study was to propose a numerical procedure which estimates the highest accumulation rate (R) of a chemical compound on the entire anatomical compartments of a mushroom body. Such data could lead to improvement in both food quality assurance and environment safety assessment. Analytical assessments on mushroom samples have shown that the accumulation potential of chemical compounds varies with mushroom species and varieties and also varies between the same mushroom anatomical compartments (see **Figure 2**) as well between mushroom development stages (“age”).

5. Solving modality path selection and motivation: the inverse numerical method for rate estimation from concentration

Inverse problems are extremely frequent in interdisciplinary science subjects. A large scale of mathematical and numerical techniques for solving scattering problems as well as other inverse problems usually exist [24]. These methods are often very different from the methods used for solving direct problems due to the differences in mathematical structure and input data [25].

In our study, the estimation process of the chemical compound accumulation rate was built on the following differential equation:

$$\frac{d(p\phi(z)\chi C(z))}{dz} - \frac{d}{dz} \left[p\phi(z) (W_c - \overline{ET}) \frac{dC(z)}{dz} \right] + p\phi(z) U_{p_f} (C(z) - B_{a_f}) = R \quad (1)$$

where the elements which may affect the rates are given as follows: z is the height [mm], ϕ is the diameter [mm], p is the porosity, χ is the hydration coefficient, W_c is the saturated hydration factor, \overline{ET} is the evaporation-transpiration coefficient, U_{p_f} is the uptake factor, B_{a_f} is the bioaccumulation factor, R is the rate of accumulation for the chemical compound of interest, C is the concentration data of the target chemical compound [$\text{ng}\cdot\text{g}^{-1}$].

The rate estimation model had as a starting point the one-dimensional transport-reaction equation for dissolved compounds presented by Lettmann et al. [26]. In this chapter the same type of equation was used but this time the equation is based on the main factors that can influence in some way the assimilation rate of a chemical compound in a vegetal organism—specifically in a mushroom body. Our model was supported by concentration data (C) of the target compounds, which were measured in the laboratory from cross-sections taken at every 2 mm over the whole body of the studied mushroom species. Also, the concentration measurements correspond to cross-sections taken at every 2 mm, and each section was divided into three concentric subintervals with regard to diameter.

The goal of the first step is to approximate R using the left-hand side of Eq. (1). The approximation of differential operators from the left side of the proposed equation has been solved using smoothing spline functions [27, 28].

Model validation was performed by solving a two-point boundary differential equation relative to Eq. (1) on the interval given by extreme values of z and comparing with the measured values of C . The numerical method for target compounds' accumulation rate validation (the solution of BVP) was implemented using MATLAB¹ `bvp4c` function [29, 30]. Good concordance was identified between the measured concentration and the concentrations computed by the solution of Eq. (1), given the rate R estimated during the validation process. The concordance is given by mean-square deviation. In the paper presented by Lettmann et al., [26] the approximation of the linear differential operator is performed by finite differences, while in our case its approximation was done through smoothing spline. Also, the rate estimation was generated randomly while in our case rate estimation was based on experimental data obtained in the laboratory. Their work, due to the nature of the practical problem, has no constraint on volume while in our case we were limited to the relative small volume and dimensions of the studied mushroom species.

6. Analysis and modeling: the estimation of R

6.1. Input data

We select for the study mushrooms from species *Macrolepiota procera*, which is one of the most popular in consumption and frequency from our country. Experimental measurements of major parameters involved in our model (parameters from Eq. (1)) were done on mature samples collected from the natural habitat where the evidence of potential contamination with chemical compounds exist. The measured data were grouped by compartments—basal bulb, stipe, gills and cap—representing the main anatomical parts of the mushroom body. Their form is illustrated in **Table 1**.

To convert our problem from bi-dimensional to a one-dimensional one, we consider a weighted mean for concentrations and parameters. The piecewise constant parameters ($p, \chi, W_c, \overline{ET}, U_{p_i}, BA_i$) are weighted by height (z), and the concentrations are weighted by diameter, on sections approximately orthogonal to median axes (see **Figure 3**).

6.2. Differential operator and its approximation

Once we have the averaged concentration we compute the rate using the formula:

$$R(z) = \frac{d(p\Phi(z)\chi C(z))}{dz} - \frac{d}{dz} \left[p\Phi(z) (W_c - \overline{ET}) \frac{dC(z)}{dz} \right] + p\Phi(z) U_{p_i} (C(z) - BA_i) \quad (2)$$

The numerical differentiations involved in Eq. (2) are critical operations, leading to large errors. Lettmann et al. [26] performed using finite differences methods followed by a Tikhonov least-squares regularization [31–33]. Our approach is different and is based on smoothing splines. The diameter ϕ is approximated by a cubic piecewise Hermite spline (MATLAB function `pchip`)

¹MATLAB is a trademark of The MathWorks, Inc.

Design parameters			
Anatomical parts	z	p	Φ
<i>Basal bulb</i>	2 - ... - 14	0.2	29 - ... - 38
<i>Stipe</i>	16 - ... - 154	0.5	41 - ... - 24
<i>Gills</i>	156 - ... - 168	0.4	160 - ... - 147
<i>Cap</i>	170 - ... - 178	0.25	143 - ... - 95
Physiological parameters			
Anatomical parts	χ	W_c	ET̄
<i>Basal bulb</i>	2.6	8.3	7.2
<i>Stipe</i>	5.2	16.6	19.2
<i>Gills</i>	7.9	38.2	28.1
<i>Cap</i>	3.6	21.7	13.5
Concentration parameters			
Anatomical parts	C	BAF	U_t
<i>Basal bulb</i>	8.1 - ... - 17.3	2	1.8
<i>Stipe</i>	18.4 - ... - 4.7	1.1	0.8
<i>Gills</i>	61.1 - ... - 69.8	2.6	1.6
<i>Cap</i>	67.9 - ... - 66.1	0.9	1.1

Table 1. Input parameters and input data structure.

and the concentration by a smoothing spline (MATLAB function *spaps*, in the spline toolbox or in the curve fitting toolbox in newer versions). Since our approximations are piecewise polynomial, the computation of their derivatives is straightforward (using *fnval* and *fnnder* functions) [34, 35]. The utilization of the smoothing spline for concentration allows us to reduce the propagated errors and to perform a correction equivalent to Tikhonov regularization [34].

6.3. The smoothing spline

We look for a spline function f , in the B-spline basis, that minimizes the expression: $\rho E_f + F(D^m f)$, where E_f is the distance of the spline function f from the given data, given by:

$$E_f = \sum_{j=1}^n w_j \|y_j - f(x_j)\|^2, \tag{3}$$

$F(D^m f)$ is:

$$F(D^m f) = \int_{x_{\min}}^{x_{\max}} \lambda(t) \|D^m f(t)\|^2 dt, \tag{4}$$

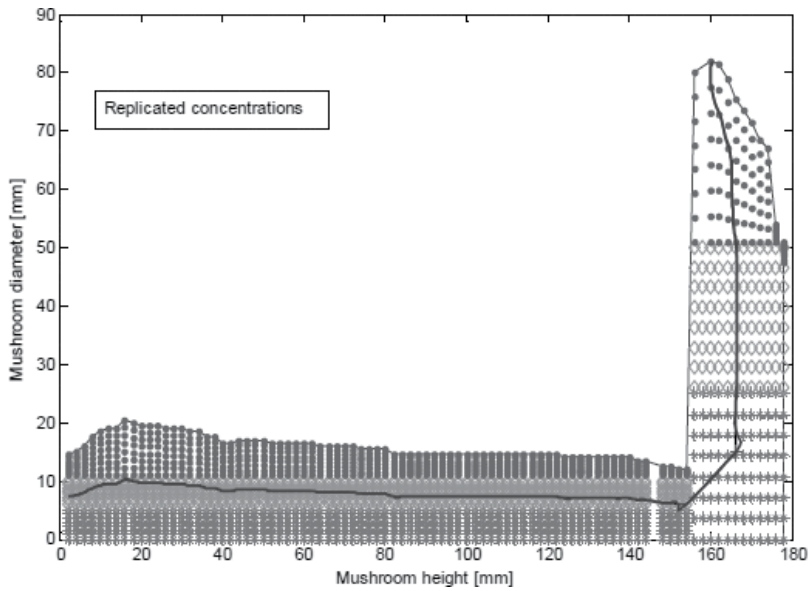


Figure 3. Schematic representation of target chemical compound concentration distribution around mushroom anatomical compartment.

and ϱ is the smoothing factor. The spaps function uses the algorithm described in Reinsch's work [32]. For additional details on smoothing and interpolation splines, see deBoor's book [33] and the MATLAB curve fitting toolbox user's guide [36, 37].

7. Validation and solution of the differential equation

This step has a double purpose: to prove that the approximation of rate, given by Eq. (2), is sufficiently accurate and to compute the concentration from the rate, without performing any measurement. We want to solve the two-point boundary value problems Eq. (1) and boundary conditions:

$$\begin{cases} C(z_{initial}) = C_0, \\ C(z_{final}) = C_n. \end{cases} \quad (5)$$

Our solution uses the collocation method [38]. The independent variable z in Eq. (1) means the length of the path along the medial axis.

8. Discussions and main conclusions

Physiological events modeling, as uptake, bioaccumulation, or metabolism, and so on, in living organisms are extremely difficult both due to the complex nature of physiological processes

and due to the complexity of the biological system that is modeled. For this reason the existence of implemented models in this area is very scarce, if not almost inexistent. Most existing models from the literature refer to models which are applied to a micro-scale fragment from a biological system (intercellular models) and less for globalized macro-scales that integrate multiple events [29, 39].

In this chapter we tried to overcome this challenge, trying to model the accumulation rate for a chemical compound based on experimental data obtained in the laboratory after analysis conducted on numerous mushroom specimen analyses of *Macrolepiota procera*. Our first approach was on the anatomical compartments of the mushroom body, but the results were not satisfactory.

Initially, the independent variable z was the height of the same compartment. Because compartments as cap and lamellae had an insufficient number of data (since the length of cap and lamellae is 1.2 cm and the minimal width for sample collection was from the section taken from 2 and 2 mm), we obtained large deviations between measured concentrations and computed concentrations at inter-compartment boundaries (see **Figure 4**).

Thus, because in chemistry in a given volume the concentration of a chemical compounds is the same in any point of these volume, we considered in the next that we have several points of concentration data of same value in horizontal sections of the cap (**Figure 3**).

These drawbacks lead us to modify the approach mentioned in Section 6 on Analysis and modeling: The estimation of R . Due to symmetry we considered a half of an axial section and a medial axis of the section. Now z is the length of the path on medial axis. To apply the classical theory on ordinary differential equations to Eq. (1), we need to have a function of class C^2 on the domain of z , while the parameters defining the rate are piecewise constant. For this reason we considered the weighted average of these parameters on the whole length of the medial axis. The diameter $\phi(z)$ was approximated by the piecewise Hermite cubic spline of measured diameters. Our choice is motivated by the fact that these splines are shape-preserving. The

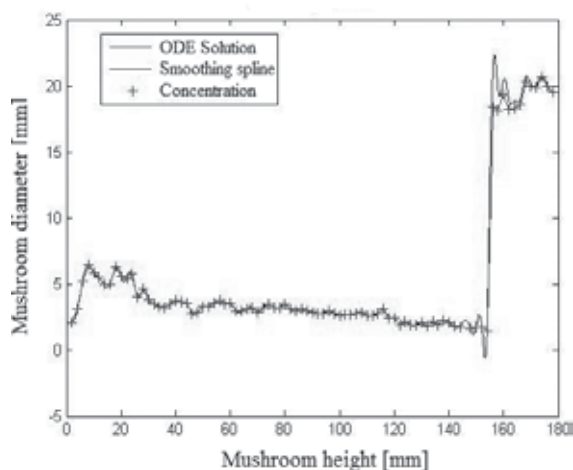


Figure 4. The averaging of parameters (ODE solution, smoothing spline) and target chemical concentrations on the whole mushroom.

next step was the application of Eq. (2) to compute the rate R . The rate R , computed along the medial axis, as described in the previous section, is plotted in **Figure 5**.

Analyzing data obtained for monitored chemical compound rate, it was possible to observe that larger fluctuations are present in the mushroom stipe while in the caps part (cap and lamellae) a decreasing tendency is registered. These data are in correlation both with infield experimental measurements and with the computed concentration obtained from our model – see **Figure 6** where the concentration C is presented, after the solution of the differential equation (1) with boundary conditions (5).

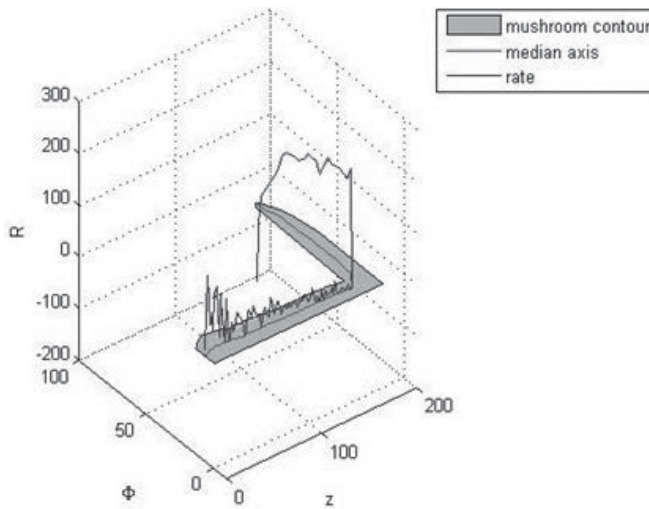


Figure 5. The graph of rate R .

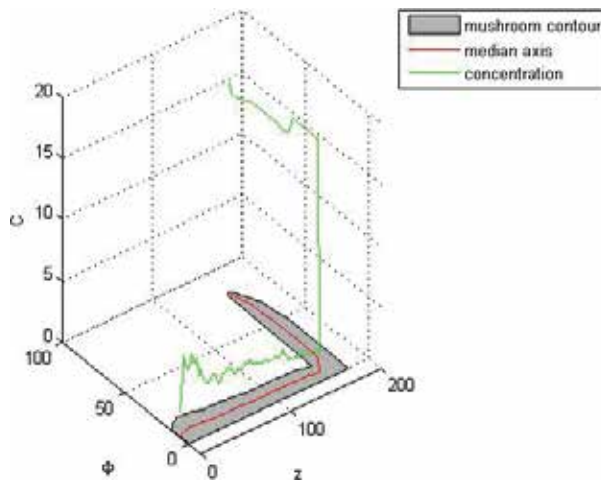


Figure 6. The graph of concentration C , obtained by the solution of a two-point boundary value problem.

In order to assess the accuracy of our model we plot the initial and the computed concentration data (**Figure 7**).

There is good matching, as **Figure 8** shows. The least square deviation for the concentration is $5.1384e-005$, and it was computed at points corresponding to the measured z and concentrations.

The model proposed for the rate estimation of a chemical compounds in living organisms, specifically a mushroom, is new. Based on our knowledge, up to date, there is no paper on the

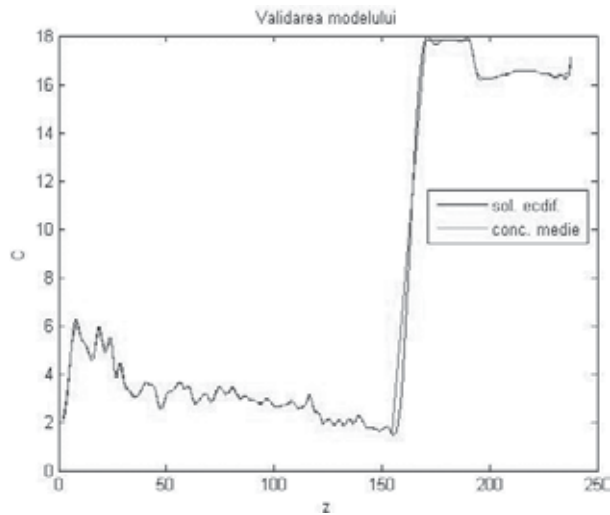


Figure 7. The graph of initial and computed concentrations.

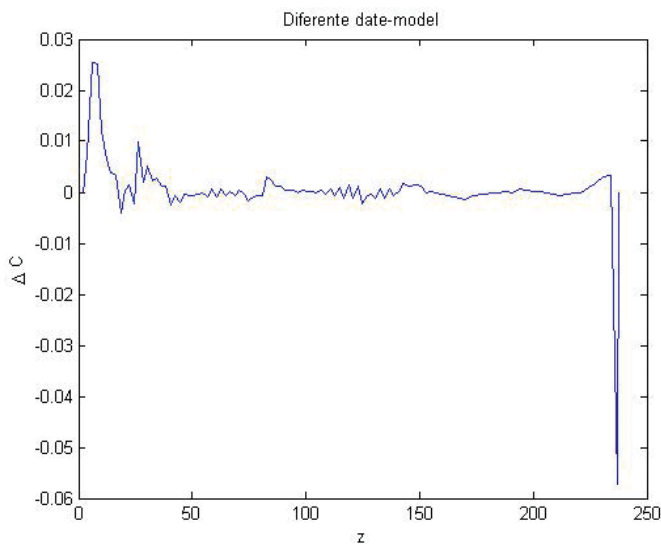


Figure 8. Differences between the initial and the computed concentration.

rate estimation for a specific chemical compound in a vegetal system. The model presented by Lettmann et al. [26] has completely different premises—there is no volume limitation in their case study while in our case we were limited to the smaller dimension of the studied living organism, the mushroom species *Macrolepiota procera*. Once have the rate for a species, we can compute concentrations via the solution of Eqs. (1) + (5) without doing any additional laboratory measurement.

Acknowledgements

This chapter was funded by the Core Program, under the support of ANCS, project OPTRONICA V number PN PN18-28 02 0.

Author details

Kovacs Eموke Dalma^{1,2*} and Kovacs Melinda Haydee²

*Address all correspondence to: dalma_kovacs2005@yahoo.com

1 Chemistry and Chemical Engineering Faculty, Babes Bolyai University, Cluj Napoca, Romania

2 National Institute for Research and Development in Optoelectronics INOE 2000, Subsidiary: Research Institute for Analytical Instrumentation, Cluj Napoca, Romania

References

- [1] Fantke P, Arnot JA, Doucette WJ. Improving plant bioaccumulation science through consistent reporting of experimental data. *Journal of Environmental Management*. 2016;**181**(1):374-384. DOI: 10.1016/j.jenvman.2016.06.065
- [2] Lebrun JD, Uher E, Tusseau-Vuillemin MH, Gourlay-Fraance C. Essential metal contents in indigenous gammarids related to exposure levels at the river basin scale: Metal-dependent models of bioaccumulation and geochemical correlation. *Science of the Total Environment*. 2014;**466-467**:100-108. DOI: 10.1016/j.scitotenv.2013.07.003
- [3] Camargo GT, Kemanian AR. Six models differ in their simulation of water uptake. *Agricultural and Forest Meteorology*. 2016;**220**:116-129. DOI: 10.1016/j.agrformet.2016.01.013
- [4] Nizzetto L, Butterfield D, Futter M, Lin Y, Allan I, Larssen T. Assessment of contaminant fate in catchments using a novel integrated hydrobiogeochemical-multimedia fate model. *Science of the Total Environment*. 2016;**544**:553-563. DOI: 10.1016/j.scitotenv.2015.11.087
- [5] Collins C. Uptake of organic pollutants and potentially toxic elements (PTEs) by crops, chapter 6. In: Martin R, Alwyn F, editors. *Persistent Organic Pollutants and Toxic Metals*

in Foods: A Volume in Series in Food Science, Technology and Nutrition. Woodhead Publishing Limited; Oxford, 2013: pp. 129-144. DOI: <https://doi-org.am.e-nformation.ro/10.1533/9780857098917.1.129>

- [6] Verkleij JAC, Golan-Goldhirsh A, Antosiewicz DM, Schwitzguebel JP, Schroder P. Dualities in plant tolerance to pollutants and their uptake and translocation to the upper plant parts. *Environmental and Experimental Botany*. 2009;**67**(1):10-22. DOI: 10.1016/j.envexpbot.2009.05.009
- [7] Siciliano SD, Germida JJ. Mechanisms of phytoremediation: Biochemical and ecological interactions between plants and bacteria. *Environmental Reviews*. 1998;**6**(1):65-79
- [8] Schroll R, Kuhn S. Test systems to establish mass balances for ¹⁴C-labeled substances in soil-plant-atmosphere systems under field conditions. *Environmental Science and Technology*. 2004;**38**(5):1537-1544
- [9] Schuelein J, Glaessgen WE, Hertkorn N, Schroeder P, sandermann H, Kettrup A. Detection and identification of the herbicide isoproturon and its metabolites in fielded samples after a heavy rainfall event. *International Journal of Environmental Analytical Chemistry*. 1996;**65**(1-4):193-202
- [10] Tynan P, Watts CD, Sowray A, Hammond I. The transport and fate of organic pollutants in rivers – II. Field measurement and modelling for styrene, xylenes, dichlorbenzenes and 4-phenyldodecane. In: Angeletti G, Bjorseth A, editors. *Organic Pollutants in Aquatic Environment*. Dordrecht: Springer; 1991. pp. 20-37
- [11] Cheng Z, Feng Y, Liu Y, Chnag H, Li Z, Xue J. Uptake and translocation of organic pollutants: A review. *Journal of Integrative Agriculture*. 2017;**16**(8):1659-1668. DOI: 10.1016/S2095-3119(16)61590-3
- [12] Brimilow RH, Chamberlain K. Principles governing uptake and transport of chemicals. In: Trapp S, Mcfarlane JC, editors. *Plant Contamination: Modelling and Simulation of Organic Chemical Processes*. Boca Raton: Lewis Publisher; 1995. pp. 38-64
- [13] Collins C, Fryer M, Grosso A. Plant uptake of non-ionic organic pollutants. *Environmental Science & Technology*. 2006;**40**(1):45-52
- [14] Liu C, Jiang X, Ma Y, Cade-Menun B. Pollutant and soil types influence effectiveness of soil-applied adsorbent in reducing rice plant uptake of persistent organic pollutants. *Pedosphere*. 2017;**27**(3):537-547. DOI: 10.1016/S1002-0160(17)60349-7
- [15] Yahyazadeh M, Nowak M, Kima H, Selmar D. Horizontal natural product transfer: A potential source of alkaloidal contaminants in phytopharmaceuticals. *Phytomedicine*. 2017;**34**:21-25. DOI: 10.1016/j.phymed.2017.07.007
- [16] Cousins IT, Mackay D. Strategies for including vegetation compartments in multimedia models. *Chemosphere*. 2001;**44**(4):643-654. DOI: 10.1016/S0045-6535(00)00514-2
- [17] Li R, Zhu Y, Zhang Y. In situ investigation of the mechanistim of the transport to tissues of polycyclic aromatic hydrocarbons adsorbed onto the root surface of *Kandelia obovata* seedlings. *Environmental Pollution*. 2015;**201**:100-106

- [18] Detternmaier EM, Doucette WJ, Bugbee B. Chemical hydrophobicity and uptake by plant roots. *Environmental Science and Technology*. 2009;**43**(2):324-329
- [19] Queguiner S, Genon LM, Roustan Y, Ciffroy P. Contribution of atmospheric emission to the contamination of leaf vegetables by persistent organic pollutants (POPs): Application to Southeastern France. *Atmospheric Environment*. 2010;**44**(7):958-967
- [20] Komp P, Mclachlan MS. Octanol/air partitioning of polychlorinated biphenyls. *Environmental Toxicology and Chemistry*. Oxford. 1997;**16**(12):2433-2437
- [21] Schroder W, Pesch R, Hertel A, Schonrock S, Harmens H, Mills G, Ilyin I. Correlation between atmospheric deposition of Cd, Hg and Pb and their concentrations in mosses specified for ecological land classes covering Europe. *Atmospheric Pollution Research*. 2013;**4**(3):267-274. DOI: 10.5094/APR.2013.029
- [22] Reis F, Barros L, Martins A, Ferreira CFR. Chemical composition and nutritional values of the most widely appreciated cultivated mushroom: An interspecies comparative study. *Food and Chemical Toxicology*. 2012;**50**:191-197
- [23] Tuo F, Zhang J, Li W, Yao S, Zhou Q, Li Z. Radionuclides in mushrooms and soil-to-mushroom transfer factors in certain areas of China. *Journal of Environmental Radioactivity*. 2017;**180**:59-64. DOI: 10.1016/j.jenvrad.2017.09.023
- [24] Kirsch A. *An Introduction to the Mathematical Theory of Inverse Problems*. 2nd ed. New York: Springer; 2011. ISBN: 978-1-4419-8473-9
- [25] Jensen DS, Wasserman A. Numerical methods for the inverse problem of density functional theory. *Quantum Chemistry*. 2017;**118**(1):1-29. DOI: 10.1002/qua.25425
- [26] Lettmann KA, Riedinger N, Ramlau R, Knab N, Bottcher ME, Khalili A, Wilff JO, Jorgensen BB. Estimation of biogeochemical rates from concentration profile: A novel inverse method. *Estuarine, Coastal and Shelf Science*. 2011. DOI: 10.1016/j.ecss.2011.01.012
- [27] Kouibia A, Pasadas M. Smoothing variational splines. *Applied Mathematics Letters*. 2000;**13**:71-75
- [28] Kano H, Fujioka H, Martin FC. Optimal smoothing and interpolating splines with constraints. *Applied Mathematics and Computation*. 2011;**218**:1831-1844
- [29] MathWorks. 2013. MATLAB Mathematics from www.mathworks.com/help/techdoc
- [30] Shampine LF, Muir PH. Estimating conditioning of BVPs for ODEs. *Mathematical and Computer Modelling*. 2004;**40**:1309-1321. DOI: 10.1016/j.mcm.2005.01.021
- [31] Wei Y, Zhang N, Ng MK, Xu W. Tikhonov regularization for weighted total least squares. *Applied Mathematics Letters*. 2007;**20**:82-87. DOI: 10.1016/j.aml.2006.03.004
- [32] Reinsch C. Smoothing by spline functions. *Numerische Mathematik*. 1967;**10**:177-183
- [33] deBoor C. *A Practical Guide to Splines*. Revised ed. New York: Springer; 2001

- [34] Lampe J, Voss H. Large-scale Tikhonov regularization of total least squares. *Journal of Computational and Applied Mathematics*. 2013;**238**:95-108. DOI: 10.1016/j.cam.2012.08.023
- [35] Ascher UM, Matheij R, Russell RD. *Numerical Solution of Boundary Value Problems for Ordinary Differential Equations*. Philadelphia: SIAM; 1995
- [36] Qu Z, Garfinkle A, Weiss NJ, Nivala M. Multi-scale modeling in biology: How to bridge the gaps between scales? *Progress in Biophysics and Molecular Biology*. 2011;**107**:21-31. DOI: 10.1016/j.pbiomolbio.2011.06.004
- [37] Fairclough HS. *Fundamentals of physiological computing. Interacting with Computers*. 2009;**21**:133-145. DOI: 10.1016/j.intcom.2008.10.011
- [38] Bica AM. Fitting data using optimal Hermite type cubic interpolating splines. *Applied Mathematics Letters*. 2012;**25**:2047-2051. DOI: 10.1016/j.aml.2012.04.016
- [39] Moler C. *Numerical Computing in MATLAB*. Philadelphia: SIAM; 2004

Applied Computational Fluid Mechanics

Benchmarks for Non-Ideal Magnetohydrodynamics

Sofronov Vasily, Zhmailo Vadim and Yanilkin Yury

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.75713>

Abstract

The paper presents an overview of benchmarks for non-ideal magnetohydrodynamics. These benchmarks include dissipative processes in the form of heat conduction, magnetic diffusion, and the Hall effect.

Keywords: ALE method, numerical simulation, magnetohydrodynamics, benchmark, verification

1. Introduction

Numerical modeling of magnetohydrodynamics (MHD) is an important and challenging problem addressed in numerous publications (e.g., see [1, 2]). This problem is further complicated in case of multi-flux models that account for the relative motion and interaction of particles of different nature (electrons, various species of ions, neutral atoms, and molecules) both with each other and with an external magnetic field.

This class of problems is generally solved using the fractional-step method, when complex operators are represented as a product of operators having a simpler structure. Thus, within the splitting method, the calculation of one-time step consists of a series of simpler procedures. It is obvious that difference schemes for each splitting stage should, where possible, preserve the properties of corresponding difference equations.

Note that the task of constructing reference solutions accounting for the whole range of physical processes is challenging (and often unfeasible). Existing benchmarks enable accuracy assessment of individual splitting stages rather than the simulation as a whole.

Magnetohydrodynamic problems are naturally divided into two groups: problems for an ideal infinitely conducting plasma and problems with dissipative processes in the form of heat conduction and magnetic viscosity.

Numerous publications on the construction of difference methods for ideal magnetohydrodynamics use a standard set of test problems. These include propagation of one-dimensional Alfvén waves at various angles to grid lines [3–5], Riemann problem for MHD equations [6–9], and various two-dimensional problems accounting for the presence of a uniform magnetic field [3, 5, 10]. In [11], a number of additional ideal MHD benchmarks are presented, which are basically shock-wave problems. A special class of tests includes problems with a weak magnetic field not affecting the medium motion. If there is an exact solution for a given hydrodynamic problem, the magnetic field “freezing-in” principle allows finding components of the field $\mathbf{H}(H_x, H_y, H_z)$ at any time with the known medium displacements $\mathbf{X} = \mathbf{X}(X_0, t)$.

The representation in publications of the problem of testing the dissipative stage of MHD equations is much the worse. Possibly, this is owing to complexity problems that require accounting the interaction of the shock-wave processes, heat conduction, diffusion of magnetic field, and Joule heating.

Numerical simulations of some of the tests presented here have been done using the Lagrangian-Eulerian code EGIDA developed at VNIIEF [12, 13] for multi-material compressible flow simulations.

The magnetohydrodynamic equation system in one-temperature approximation modified by the Hall effect can be written in the following conservative form [2]:

$$\begin{aligned} \frac{\partial \rho}{\partial t} + \operatorname{div} \rho \mathbf{u} &= 0, & \frac{\partial \rho \mathbf{u}}{\partial t} + \operatorname{div}(\rho \mathbf{u} \otimes \mathbf{u} + (P + P_H) \mathbf{I} - 0.5 \mathbf{H} \otimes \mathbf{H}) &= 0, & P_H &= 0.5 |\mathbf{H}|^2, \\ \frac{\partial \mathbf{H}}{\partial t} + \operatorname{div}(\mathbf{u} \otimes \mathbf{H} - \mathbf{H} \otimes \mathbf{u}) &= -\operatorname{rot}(v \cdot \operatorname{rot} \mathbf{H} + b[\mathbf{H} \times \operatorname{rot} \mathbf{H}]), & \frac{\partial \rho n_e}{\partial t} + \operatorname{div}(\rho n_e \mathbf{u}) &= 0, & (1) \\ \frac{\partial \Xi}{\partial t} + \operatorname{div}((\Xi + P + P_H) \mathbf{u} - \mathbf{H}(\mathbf{u} \cdot \mathbf{H}) - \kappa \operatorname{grad} T) &= 0, & \Xi &= \rho(e + 0.5 |\mathbf{u}|^2) + P_H, \\ P &= P(\rho, T), & \varepsilon &= \varepsilon(\rho, T). \end{aligned}$$

where $v = c^2/(4\pi\sigma)$ is the magnetic viscosity coefficient, κ is the heat conduction factor, $b = c/(4\pi en_e)$ is a local exchange (Hall) parameter [2], and e and n_e are charge and density of electrons. When writing Eq. (18), we assume that bias currents and electron inertia are negligibly small [2]. Equation system (1) differs from equation system for ideal MHD owing to diffusion terms present in the equations of energy and inductance of magnetic field.

2. A plane diffusion wave with regard to the Hall effect

Let the components of magnetic field depend on coordinate z only, i.e., $\mathbf{H} = (H_x(t, z), H_y(t, z), H_{z0})$. We neglect the medium motion. Then, the magnetic field equation (for components) is written in the following form:

$$\frac{dH_x}{dt} = \nu \frac{\partial^2 H_x}{\partial z^2} + \beta \frac{\partial^2 H_y}{\partial z^2}, \quad \frac{dH_y}{dt} = \nu \frac{\partial^2 H_y}{\partial z^2} - \beta \frac{\partial^2 H_x}{\partial z^2}, \quad \frac{dH_z}{dt} = 0, \quad \beta = bH_{z0}. \quad (2)$$

Consider the problem of a diffusion wave propagating in an unbounded medium with the given boundary and initial conditions:

$$\mathbf{H}(t, z = 0) = \mathbf{H}_1, \quad \mathbf{H}(t, z \rightarrow \infty) = \mathbf{H}_0, \quad \mathbf{H}(t = 0, z) = \mathbf{H}_0, \quad \mathbf{H}_0 = (0, 0, H_{z0}), \quad \mathbf{H}_1 = (H_{x0}, H_{y0}, H_{z0}). \quad (3)$$

Let $\gamma = \sqrt{\nu^2 + \beta^2}$. A general solution to Eq. (2) for the self-similar variable $\xi = z/\sqrt{4\gamma t}$ looks like

$$H_x = H_{x0} + C_1\Phi(\xi) + C_2\Psi(\xi), \quad H_y = H_{y0} + C_1\Psi(\xi) - C_2\Phi(\xi), \quad H_z = H_{z0}$$

where $\Phi(\xi) = \int_0^\xi \exp(-\nu x^2/\gamma) \sin(\beta x^2/\gamma) dx$, $\Psi(\xi) = \int_0^\xi \exp(-\nu x^2/\gamma) \cos(\beta x^2/\gamma) dx$.

Since $\Phi(\infty) = \Gamma \frac{1}{2} \sqrt{\frac{\gamma-\nu}{2\gamma}} = \sqrt{\frac{\pi(\gamma-\nu)}{2\gamma}}$, $\Psi(\infty) = \Gamma \frac{1}{2} \sqrt{\frac{\gamma+\nu}{2\gamma}} = \sqrt{\frac{\pi(\gamma+\nu)}{2\gamma}}$ constants C_1, C_2 with regard to boundary conditions can be found from equations

$$C_1 = \frac{-H_{x0}\Phi(\infty) - H_{y0}\Psi(\infty)}{\Phi(\infty) + \Psi(\infty)} = -\frac{2}{\sqrt{\pi}} \left(H_{x0} \sqrt{\frac{\gamma-\nu}{2\gamma}} + H_{y0} \sqrt{\frac{\gamma+\nu}{2\gamma}} \right),$$

$$C_2 = \frac{-H_{x0}\Psi(\infty) + H_{y0}\Phi(\infty)}{\Phi(\infty) + \Psi(\infty)} = -\frac{2}{\sqrt{\pi}} \left(H_{x0} \sqrt{\frac{\gamma+\nu}{2\gamma}} - H_{y0} \sqrt{\frac{\gamma-\nu}{2\gamma}} \right).$$

Simulation setup: the initial data is described by Eq. (3). A bounded computational domain $0 < z < L, L = 1$ is considered. For this reason, the magnetic field value taken from the analytical solution $H_x(t, z = L) = 1 + (C_1\Phi + C_2\Psi) \frac{L}{\sqrt{4\gamma t}}$, $H_y(t, z = L) = 1 + (C_1\Psi - C_2\Phi) \frac{L}{\sqrt{4\gamma t}}$, $H_z(t, z =$

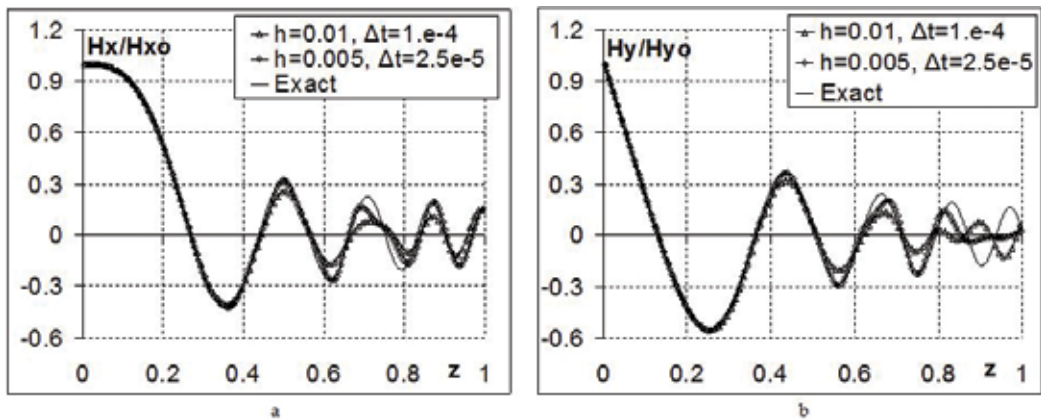


Figure 1. Profiles of field components at time $t = 0.01$: (a) H_x and (b) H_y .

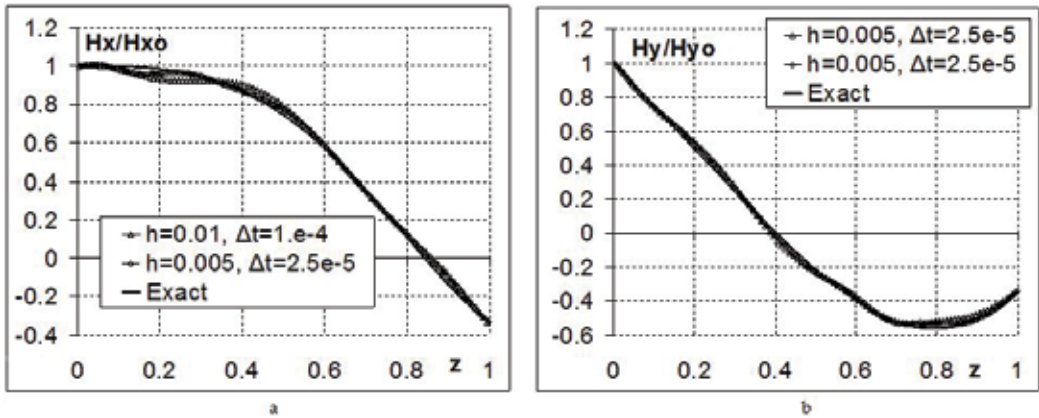


Figure 2. Profiles of field components at time $t = 0.1$: (a) H_x and (b) H_y .

$L) = 1$ is imposed on the right boundary $z = L$. On the left boundary $z = 0$, the field components take constant values according to Eq. (3). In simulations with 2D and 3D codes, boundary conditions $\partial \mathbf{H} / \partial \mathbf{n} = 0$ (\mathbf{n} is a normal vector to a face) are imposed on lateral faces. By varying parameters ν and β , we can study the effect of the diffusion and Hall terms in Eq. (2) on the diffusion wave parameters. Consider an option with the Hall effect dominating over the diffusion effect: $\nu = 0$, $\beta = 1$, $H_{x0} = H_{y0} = H_{z0} = 1$. Profiles of magnetic field's components H_x , H_y at time $t = 0.01, 0.1$ are shown in Figures 1 and 2. With the grid refinement, convergence to the reference solution takes place.

3. Diffusion of magnetic field in an immovable plane layer of plasma with regard to joule heating and its effect on the diffusion and heat conduction coefficients

The problem of magnetic diffusion in a plane layer of material has many applications in practice [14]. In its detailed formulation, the problem was considered in paper [14] for mega gauss fields. Hydrodynamic motion, magnetic diffusion, heat conduction by electrons, and radiant heat exchange in the "back and forth" approximation were taken into account. Since finding an exact solution to such a problem causes difficulties, the original formulation needs to be simplified. Self-similar solutions to the problem obtained with simplifying assumptions were also presented in [15].

A model problem is considered with the following assumptions:

- plasma is immovable, it has a constant heat capacity,
- plasma has Coulomb conductivity,
- heat conduction is absent.

With such assumptions, the problem is reduced to solving equations

$$\frac{d\mathbf{H}}{dt} = -\text{rot}(\nu \cdot \text{rot}\mathbf{H}), \quad \rho \frac{dT}{dt} = (\gamma - 1)\nu(\text{rot}\mathbf{H} \cdot \text{rot}\mathbf{H}), \quad (4)$$

where $\nu = c^2/4\pi\sigma$, $\sigma = \sigma_0(T/Ry)^\alpha$, $\alpha = 3/2$, $\gamma - 1 = R/C_V$, $R = 1$, $C_V = 1.5$

Here, ρ is the density of plasma, γ is the heat capacity ratio (adiabatic index), σ_0 is the conductivity at $T = Ry$ (it is expressed via atomic constants), and Ry is the Rydberg constant. Energy units have been chosen for temperature.

At initial time $t = 0$, all quantities depend on one space coordinate. It is assumed that the magnetic field has only one component, $\mathbf{H} = (0, 0, H_z)$. The solution is considered for the problem with initial conditions Eq. (5) and boundary conditions Eq. (6):

$$H_z(x, t = 0) = \begin{cases} 0 & \text{if } x < 0 \\ H_0 & \text{if } x > 0 \end{cases}, \quad T(x, t = 0) = T_0, \quad \rho(x, t = 0) = \rho_0 \quad (5)$$

$$H_z(x \rightarrow -\infty, t) = 0, \quad H_z(x \rightarrow \infty, t) = H_0, \quad T(x \rightarrow \pm\infty, t) = T_0. \quad (6)$$

For dimensionless variables, $h_z = H_z/H_0$, $\tau = T/T_0$ Eq. (4) are reduced to the form:

$$\frac{dh_z}{dt} = \frac{\partial}{\partial x} \nu(\tau) \frac{\partial h_z}{\partial x}, \quad \frac{d\tau}{dt} = \eta \nu(\tau) \left(\frac{\partial h_z}{\partial x} \right)^2, \quad \nu(\tau) = \nu_0 \tau^{-\alpha}, \quad \nu_0 = \frac{c^2}{4\pi\sigma_0} \left(\frac{T_0}{Ry} \right)^\alpha, \quad \eta = (\gamma - 1) \frac{H_0^2}{\rho T_0}. \quad (7)$$

In an infinite region ($-\infty < x < \infty$), the problem has a self-similar solution depending on the variable $\xi = x/\sqrt{\nu_0 t}$. The solution can be obtained by integrating the system of ordinary differential equations:

$$\frac{\xi}{2} \frac{dh_z}{d\xi} + \frac{d}{d\xi} \left(\tau^{-\alpha} \frac{dh_z}{d\xi} \right) = 0, \quad \frac{\xi}{2} \frac{d\tau}{d\xi} + \eta \tau^{-\alpha} \left(\frac{dh_z}{d\xi} \right)^2 = 0. \quad (8)$$

with boundary conditions

$$h_z(\xi \rightarrow -\infty) = 0, \quad h_z(\xi \rightarrow \infty) = 1, \quad \tau(\xi \rightarrow \pm\infty) = 1. \quad (9)$$

Note that for a linear case, $\alpha = 0$, the solution of Eqs. (8), (9) can be found in quadratures

$$h_z(\xi) = 0.5(1 + \text{sign}(\xi)\text{erf}(\xi/2)), \quad \tau(\xi) = 1 - \frac{\eta}{4\pi} Ei(-\xi^2/4). \quad (10)$$

Since $Ei(-x) = C + \ln x + \sum_i \frac{(-1)^i x^i}{i!}$, temperature in the vicinity of interface $\xi^2 \sim 0$ for the linear case $\alpha = 0$ has the logarithmic profile $\tau(\xi) \sim -\eta \ln \xi^2/4\pi$.

In general, if $\alpha > 0$, one does not manage to establish the asymptotic law for temperature in the vicinity of $\xi^2 \sim 0$, because of no integral curves of Eq. (8) satisfying the boundary conditions at infinity Eq. (9).

Now, let us build the reference solution to the problem with regard to heat conduction. In this case temperature near the interface takes a finite value. The diffusion equations and energy equation of magnetic field with regard to Joule heating and heat conduction are considered. As it was assumed earlier, all quantities depend on one space coordinate, and the magnetic field has only one component, $\mathbf{H} = (0, 0, H_z)$. For dimensionless variables, $h_z = H_z/H_0$ and $\tau = T/T_0$ equations are reduced to the forms

$$\frac{dh_z}{dt} = \frac{\partial}{\partial x} \nu_0 \tau^{-\alpha} \frac{\partial h_z}{\partial x}, \quad \frac{d\tau}{dt} = \eta \nu_0 \tau^{-\alpha} \left(\frac{\partial h_z}{\partial x} \right)^2 + \frac{\partial}{\partial x} \kappa_0 \tau^\beta \frac{\partial \tau}{\partial x}. \quad (11)$$

A self-similar solution depending on the variable $\xi = x/\sqrt{\nu_0 t}$ can be obtained by integrating the system of ordinary differential equations:

$$\frac{\xi}{2} \frac{dh_z}{d\xi} + \frac{d}{d\xi} \left(\tau^{-\alpha} \frac{dh_z}{d\xi} \right) = 0, \quad \frac{\xi}{2} \frac{d\tau}{d\xi} + \eta \tau^{-\alpha} \left(\frac{dh_z}{d\xi} \right)^2 + a \frac{d}{d\xi} \tau^\beta \frac{d\tau}{d\xi} = 0, \quad a = \frac{\kappa_0}{D_0}, \quad (12)$$

with boundary conditions:

$$h_z(\xi \rightarrow -\infty) = 0, \quad h_z(\xi \rightarrow \infty) = 1, \quad \tau(\xi \rightarrow \pm\infty) = 1. \quad (13)$$

To find the reference solution, it is convenient to use the first-order system with an increased number of unknowns instead of the second-order system Eq. (12). The first-order system relative to variables $h_z, \tau, \Psi = \tau^{-\alpha} dh_z/d\xi, w = -a\tau^\beta d\tau/d\xi$ looks like

$$\frac{dh_z}{d\xi} = \Psi \tau^\alpha, \quad \frac{d\tau}{d\xi} = -\frac{w\tau^{-\beta}}{a}, \quad \frac{d\Psi}{d\xi} = -\frac{\xi\Psi\tau^\alpha}{2}, \quad \frac{dw}{d\xi} = -\frac{\xi w\tau^{-\beta}}{2a} + \eta\Psi^2\tau^\alpha. \quad (14)$$

Consider the numerical solution of Eq. (14) for the right half plane ($0 < \xi < \infty$). The solution in the left half plane ($-\infty < \xi < 0$) follows from the symmetry conditions:

$$h_z(-\xi) = 1 - h_z(\xi), \quad \tau(-\xi) = \tau(\xi), \quad \Psi(-\xi) = \Psi(\xi), \quad w(-\xi) = -w(\xi).$$

Consider the numerical solution of Eq. (14) in a bounded domain ($0 < \xi < \xi_1$). To formulate boundary conditions for this bounded domain, it is required to find the asymptotic behavior of functions with $\xi \rightarrow \infty$. Asymptotic laws can be formulated, if we assume $a = 0.5$. In this case boundary conditions have the form:

$$h_z(\xi \rightarrow \infty) = 0.5(1 + \operatorname{erf}(\xi/2)), \quad \tau(\xi \rightarrow \infty) = 1, \quad (15)$$

$$\Psi(\xi \rightarrow \infty) = c_1 \exp(-\xi^2/4), \quad w(\xi \rightarrow \infty) = (c_2 + \eta\xi c_1^2) \exp(-\xi^2/2).$$

Constants C_1, C_2 are taken so that the following conditions are satisfied on the left boundary of the computational domain:

$$h_z(\xi = 0) = 0.5, \quad w(\xi = 0) = 0.$$

Confine oneself to the consideration of case $a = 0.5$. Introduce new variables $\Psi(\xi) = \Psi(\xi) \exp(\xi^2/4)$, $W(\xi) = w(\xi) \exp(-\xi^2/2)$ with regard to boundary conditions Eq. (15). The replacement of variables gives us the equation system:

$$\begin{aligned} \frac{dh_z}{d\xi} &= \Psi \tau^\alpha \exp\left(-\frac{\xi^2}{4}\right), & \frac{d\tau}{d\xi} &= -\frac{W \tau^{-\beta}}{a} \exp\left(-\frac{\xi^2}{2}\right), \\ \frac{d\Psi}{d\xi} &= -\frac{\xi \Psi (1 - \tau^\alpha)}{2}, & \frac{dW}{d\xi} &= \frac{\xi W (1 - \tau^{-\beta})}{2a} + \eta \Psi^2 \tau^\alpha \end{aligned} \quad (16)$$

with boundary conditions:

$$\begin{aligned} h_z(\xi = \xi_1) &= 0.5(1 + \operatorname{erf}(\xi_1/2)), & \tau(\xi = \xi_1) &= 1, \\ \Psi(\xi = \xi_1) &= c_1, & W(\xi = \xi_1) &= c_2 + \eta \xi_1 c_1^2, \\ h_z(\xi = 0) &= 0.5, & w(\xi = 0) &= 0. \end{aligned} \quad (17)$$

The set of Eqs. (16), (17) was solved numerically with the methods of automatically selecting an integration step. The following values of parameters were used in simulations: $\xi_1 = 10$, $\eta = 20/3$, $\alpha = 3/2$, $D_0 = 1$, and $k_0 = aD_0 = 1/2$. The values of constants satisfying the boundary conditions Eq. (17) were obtained: $C_1 = 0.10231$ and $C_2 = 1.79474$. Since the behavior of functions near the right boundary corresponds to asymptotic laws Eq. (17), the reduction of parameter ξ_1 from $\xi_1 = 10$ to $\xi_1 = 1$ does not affect simulation results.

Results of simulations are illustrated in **Figures 3** and **4**. With the use of such regularity method (with regard to heat conduction), temperature at the central point of the computational domain takes its finite value. Note that with $t = 1/\nu_0$ the space coordinate coincides with the self-similar coordinate, $x = \xi$.

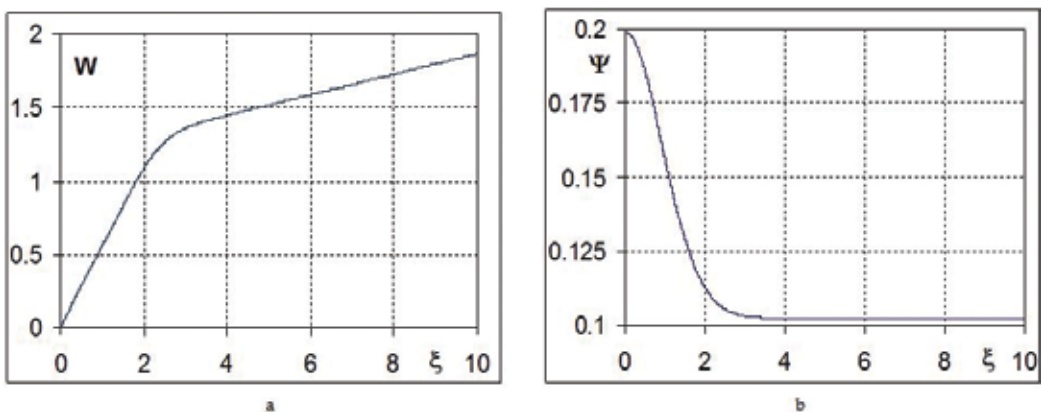


Figure 3. Profiles of self-similar functions: (a) W and (b) Ψ .

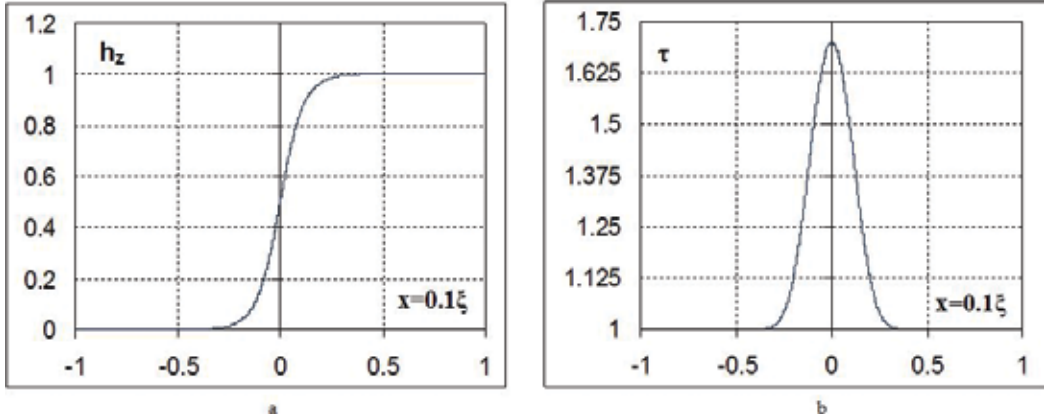


Figure 4. Profiles: (a) magnetic field and (b) temperature.

4. A point explosion in a perfectly non-conducting atmosphere

Let us consider the problem of a point blast in the presence of a uniform magnetic field (for definiteness) along the z axis ($H_z = H_{z0} = 0.01$) in a perfectly non-conducting atmosphere. Initial data are chosen in such a way that the field has no effect on the motion of matter $\varepsilon_0(r_2/r_1)^3 \gg H_{z0}^2/\rho_0$. It is assumed that behind the shock front, the medium becomes perfectly conducting. If a self-similar solution to the problem of a point explosion is known, then one can calculate the magnetic field components at some time $t > 0$. In an external domain ($r > r_F(t)$), the magnetic field's vector potential is the solution to stationary equation $rot(rot\Psi) = 0$. With regard to conditions at infinity, this solution takes the form [16]:

$$\Psi_r(r, \vartheta, t) = 0, \quad \Psi_\varphi(r, \vartheta, t) = 0.5H_{z0}r(1 - C(t)/r^3) \sin \vartheta, \quad \Psi_\vartheta(r, \vartheta, t) = 0. \quad (18)$$

Here, unknown constant $C(t) = br_F^3(t)$ can be found from the condition of coupling with the solution in an internal domain ($r < r_F(t)$). Write components of magnetic field $\mathbf{H} = rot\Psi$:

$$H_r(r, \vartheta, t) = H_{z0}(1 - C(t)/r^3) \sin \vartheta, \quad H_\varphi(r, \vartheta, t) = 0, \quad H_\vartheta(r, \vartheta, t) = -H_{z0}(1 + 0.5C(t)/r^3) \cos \vartheta.$$

The solution in the internal domain ($r < r_F(t)$) is found from the freezing-in condition of the magnetic field:

$$\frac{d}{dt} \left(\frac{H_r}{\rho} \right) = \frac{H_r}{\rho} \frac{\partial u_r}{\partial r}, \quad \frac{d}{dt} \left(\frac{H_\vartheta}{\rho} \right) = \frac{H_\vartheta}{\rho} \frac{u_r}{r}, \quad \frac{d}{dt} \left(\frac{H_\varphi}{\rho} \right) = \frac{H_\varphi}{\rho} \frac{u_r}{r}.$$

The integration of these equations with regard to the solution in external domain Eq. (18) gives us

$$H_r(r, \vartheta, t) = H_{z0}h_r(r, t) \cos \vartheta, \quad H_\vartheta(r, \vartheta, t) = H_{z0}h_\vartheta(r, t) \sin \vartheta, \quad H_\varphi(r, \vartheta, t) = 0, \quad (19)$$

where

$$h_r(r, t) = \begin{cases} 1 - \beta \left(\frac{r_F(t)}{r}\right)^3, & r > r_F(t) \\ (1 - \beta) \left(\frac{r_0(r, t)}{r}\right)^2, & r \leq r_F(t) \end{cases}, h_\vartheta(r, t) = \begin{cases} -1 - \frac{\beta}{2} \left(\frac{r_F(t)}{r}\right)^3, & r > r_F(t) \\ -\left(1 + \frac{\beta}{2}\right) \frac{\rho(r, t)r(\gamma - 1)}{\rho_0 r_0(r, t)(\gamma + 1)}, & r \leq r_F(t) \end{cases}.$$

Here, the functions $r_0(r, t)$, $\rho(r, t)$ are defined from the self-similar solution to the point blast problem [17].

Unknown constant β can be found from the condition of the solenoidal distribution of magnetic field in internal domain.

$$\text{div}\mathbf{H} = H_{z0} \cos \vartheta \left(\frac{1}{r^2} \frac{\partial r^2 h_r}{\partial r} + 2 \frac{h_\vartheta}{r} \right) = 2 \frac{H_{z0}}{r_0} \cos \vartheta \left((1 - \beta) \frac{r_0^2 \partial r_0}{r^2 \partial r} - \left(1 + \frac{\beta}{2}\right) \frac{\rho}{\rho_0} \frac{\gamma - 1}{\gamma + 1} \right) = 0.$$

Since $\frac{\rho}{\rho_0} = \frac{r_0^2 \partial r_0}{r^2 \partial r}$, then $1 - \beta = \left(1 + \frac{\beta}{2}\right) \frac{\gamma - 1}{\gamma + 1}$. That is why $\beta = \frac{4}{3\gamma + 1}$.

Note that in external domain this condition is satisfied automatically. In Cartesian coordinates, the solution of Eq. (19) looks like

$$\begin{aligned} H_x(x, y, z, t) &= \frac{H_{z0}xz}{r^2} (h_r(r, t) + h_\vartheta(r, t)), & H_y(x, y, z, t) &= \frac{H_{z0}yz}{r^2} (h_r(r, t) + h_\vartheta(r, t)), \\ H_z(x, y, z, t) &= H_{z0} \left(\frac{z^2}{r^2} (h_r(r, t) + h_\vartheta(r, t)) - h_\vartheta(r, t) \right). \end{aligned} \tag{20}$$

It is convenient to compare the numerical and exact solutions using the field components depending on one space coordinate:

$$\begin{aligned} h_r(r, t) &= \frac{H_r}{H_{z0} \cos \vartheta} = \frac{1}{H_{z0}} \left(H_z(x, y, z, t) + \frac{xH_x(x, y, z, t) + yH_y(x, y, z, t)}{z} \right), \\ h_\vartheta(r, t) &= \frac{H_\vartheta}{H_{z0} \sin \vartheta} = -\frac{1}{H_{z0}} \left(H_z(x, y, z, t) - \frac{z(xH_x(x, y, z, t) + yH_y(x, y, z, t))}{r^2 - z^2} \right), \\ h_\varphi(r, t) &= \frac{H_\varphi}{H_{z0} \sin \vartheta} = \frac{1}{H_{z0}} \left(\frac{-yH_x(x, y, z, t) + xH_y(x, y, z, t)}{r} \right) = 0. \end{aligned} \tag{21}$$

The magnetic field lines can be obtained by integrating equations

$$\frac{dx}{dz} = \frac{xz(h_r(r, t) + h_\vartheta(r, t))}{-r^2 h_\vartheta(r, t) + z^2(h_r(r, t) + h_\vartheta(r, t))}, \quad \frac{dy}{dz} = \frac{yz(h_r(r, t) + h_\vartheta(r, t))}{-r^2 h_\vartheta(r, t) + z^2(h_r(r, t) + h_\vartheta(r, t))}. \tag{22}$$

We consider the process stage, at which the numerical simulation becomes self-similar. In this case, the shock wave is considerably far (compared to the energy release region) from the blast center. For example, at the final time $t = 3$, the wave front is located at a distance of $R_F = 13.467$.

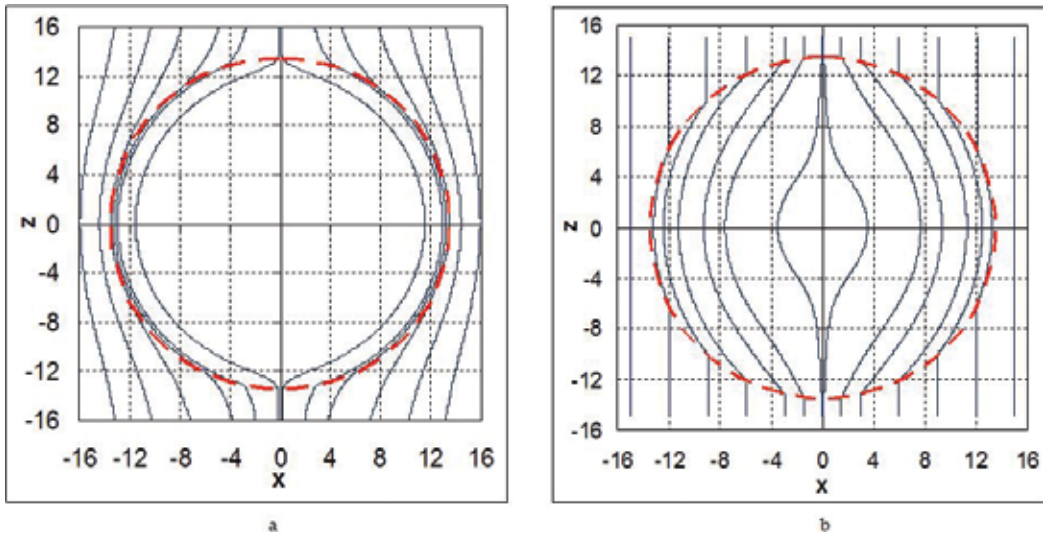


Figure 5. The magnetic field lines at time $t = 3$ in plane $y = 0$. A dashed line shows the shock front. Strong blast in a perfectly nonconducting atmosphere (a) and in a uniform conducting atmosphere (b) [11].

The flow parameters in this problem depend on a single spatial variable, r , and the field components, on two variables, r and θ . One can restrict the consideration to any plane passing through the z axis. Magnetic field lines in the plane $y = 0$ at $t = 3$ are shown in **Figure 5**. It follows from this figure that these lines of force stretch along axis z and, thereby, prevent the spread of gas in the direction orthogonal to this axis. This effect is small in the given problem because of the field smallness. With an increased strength of the field, the pressure zone gets out of its spherical shape due to occurrence of the singled out direction.

Profiles of the field components H_y and H_z along the line $x = z$ in this plane are presented in **Figure 6**. The self-similar functions of fluid parameters and radial and angular field components h_r and h_θ depend on one spatial variable r . These functions are shown in **Figure 7**. In testing numerical methods, the values of grid functions for all cells in the domain can be mapped onto the

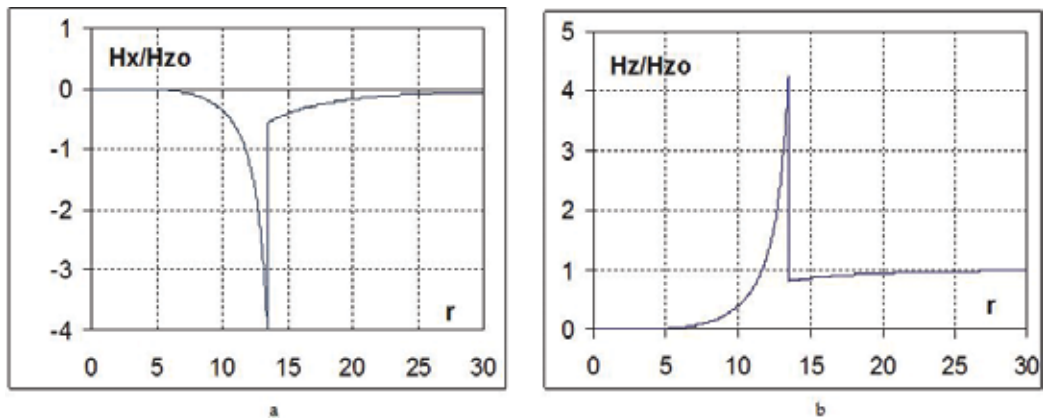


Figure 6. Profiles of field components along line $x = z$ and $y = 0$ at time $t = 3$: (a) H_x and (b) H_z .

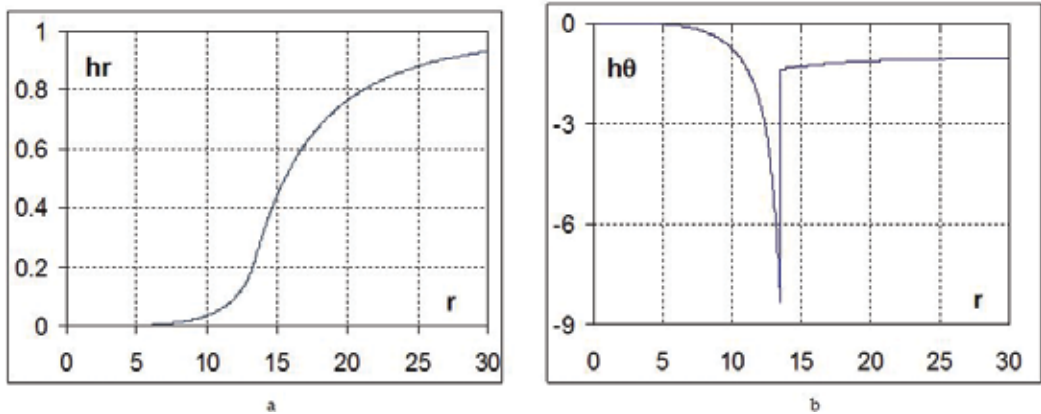


Figure 7. Profiles of the non-dimensionalized field components at time $t = 3$: (a) radial h_r , and (b) angular h_θ .

figures. Such a comparison of reference versus numerical solution indicates whether the spherical symmetry is preserved during numerical simulations.

Simulation setup: The energy release region is a sphere of radius $r_1 = 0.1$, in which the initial specific energy per unit mass is set to $\varepsilon = \varepsilon_0 = 10^7$ and the initial density is set to $\rho = \rho_0 = 1$. In the spherical layer $r_1 < R < r_2 = 15$, the specific energy and the density are equal to $\varepsilon = 0, \rho = \rho_0 = 1$, respectively. The EOS is $P = \rho(\gamma - 1)\varepsilon, \gamma = 1.4$. The problem domain in the three-dimensional setup is a cube $L \times L \times L$. All boundary faces of the domain are rigid walls.

This problem requires taking into account the magnetic field diffusion in external domain (outside the shock front). It is assumed that behind the shock front, the medium becomes perfectly conducting due to ionization effects. The magnetic viscosity is approximated by the following dependence:

$$v(\varepsilon) = \frac{c^2}{4\pi\sigma} = \begin{cases} v_1 = 10^5, & \varepsilon \leq 0, \\ v_1(1 - \varepsilon/\varepsilon_1), & 0 < \varepsilon < \varepsilon_1 = 1, \\ 0, & \varepsilon_1 \leq \varepsilon. \end{cases}$$

Parameter ε_1 is chosen to provide that the magnetic viscosity behind the shock front is always zero, i.e., the condition $\varepsilon_1 < \varepsilon_F(t) \simeq \varepsilon_0(r_1/r_F(t))^3, t \leq t_k = 3$ is satisfied. Accounting of diffusion in external domain leads to the necessity of increasing the size of computational domain $(r_F(t_k)/L)^3 < < 1$) in comparison with the ideal MHD problems to be able to set boundary conditions corresponding to the initial undisturbed state. Thus, the size of computational domain must be almost five times larger, $L \approx 75$.

5. Diffusion of magnetic field into a spherical plasma cloud

The problem formulation and its analytical solution have been taken from [16]. In contrast to this paper, consider the diffusion problem (the plasma cloud motion is neglected):

$$\frac{\partial \mathbf{H}}{\partial t} = -\text{rot}(v \cdot \text{rot} \mathbf{H} + b[\mathbf{H} \times \text{rot} \mathbf{H}]). \quad (23)$$

It is assumed that the magnetic field at infinity is uniform and directed along axis z : $\mathbf{H}(0, 0, H_0 = \sqrt{2})$ (see **Figure 8**). The magnetic viscosity coefficient is constant inside and outside the cloud: $v(r) = \begin{cases} v_1, & r < r_0 = 1 \\ v_2, & r > r_0 \end{cases}$.

5.1. Diffusion of magnetic field in the absence of the Hall effect

Assume that the Hall effect contribution is small, $bH_0/v \ll 1$. Write the equation of diffusion relative to vector potential $\mathbf{H} = \text{rot} \Psi$:

$$\frac{\partial \Psi}{\partial t} = -v \cdot \text{rot} \text{rot} \Psi. \quad (24)$$

This is an axially symmetric problem, and, therefore, it is convenient to use the polar coordinate system (r, ϑ, φ) . With no Hall effect and with regard to the conditions at infinity, the initial data for vector potential takes the form [16]

$$\Psi_r(r, \vartheta, t = 0) = 0, \quad \Psi_\vartheta(r, \vartheta, t = 0) = 0, \quad \Psi_\varphi(r, \vartheta, t = 0) = rf(r, t = 0) \sin \vartheta, \quad (25)$$

$$f(r, t = 0) = \frac{H_0}{2} \begin{cases} 0 & r < r_0 \\ (1 - (r_0/r)^3) & r > r_0 \end{cases}. \quad (26)$$

The solution of Eq. (24) with initial data Eq. (25) is reduced to solving Eq. (27) with initial data Eq. (26) and boundary conditions Eq. (28):

$$\frac{df}{dt} = \frac{v(r)}{r^4} \frac{\partial}{\partial r} r^4 \frac{\partial f}{\partial r}, \quad (27)$$

$$\frac{\partial f}{\partial r}(r = 0, t) = 0, \quad \frac{\partial f}{\partial r}(r \rightarrow \infty, t) = 0. \quad (28)$$

The solution has the form

$$\Psi_r(r, \vartheta, t = 0) = 0, \quad \Psi_\vartheta(r, \vartheta, t) = 0, \quad \Psi_\varphi(r, \vartheta, t) = rf(r, t) \sin \vartheta. \quad (29)$$

Paper [16] considers the plasma cloud interaction with the magnetic field of vacuum, and, therefore, $v_2 \rightarrow \infty$ is assumed. For this special case, the solution to Eq. (27) in quadratures has been obtained, and it has the forms:

$$f(r, t) = \frac{H_0}{2} \begin{cases} 1 - \frac{6}{\zeta^2} \sum_{n=1}^{\infty} \frac{(-1)^n}{(\pi n)^2} T_n(t) \cdot \left(\cos \pi n \zeta - \frac{\sin \pi n \zeta}{\pi n \zeta} \right), & 0 < \zeta = \frac{r}{r_0} < 1 \\ 1 - \frac{6}{\zeta^2} \sum_{n=1}^{\infty} \frac{1}{(\pi n)^2} T_n(t), & 1 < \zeta \end{cases}, \quad T_n(t) = \exp\left(-(\pi n)^2 \frac{v_1 t}{r_0^2}\right). \quad (30)$$

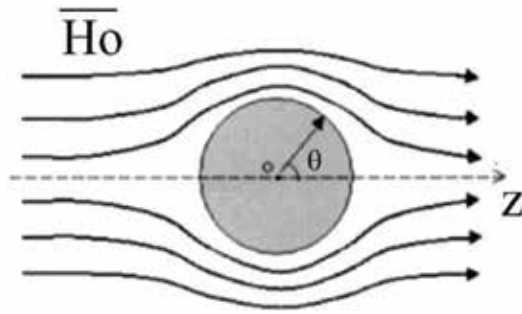


Figure 8. The problem of diffusion into a plasma cloud [16].

For finite values of conductivity in external domain $\nu_2 > 0$, the limit numerical solution of Eq. (27) has been taken for the reference solution.

Components of magnetic field $\mathbf{H} = \text{rot}\Psi$ are found by differentiating vector potential Eq. (29). If function $f(r, t)$ is known, these components are calculated using formulas.

$$\begin{aligned} H_r(r, \vartheta, t) &= h_r(r, t) \cos \vartheta, H_\vartheta(r, \vartheta, t) = -h_\vartheta(r, t) \sin \vartheta, H_\varphi(r, \vartheta, t) = h_\varphi \sin \vartheta, h_r(r, t) \\ &= 2f(r, t), h_\vartheta(r, t) = \partial r^2 f / r \partial r, h_\varphi(r, t) = 0 \end{aligned} \quad (31)$$

In Cartesian coordinates the field components have the forms

$$\begin{aligned} H_z(x, y, z, t) &= h_\vartheta(r, t) + \frac{z^2}{r^2} (h_r(r, t) - h_\vartheta(r, t)), H_y(x, y, z, t) \\ &= \frac{zy}{r^2} (h_r(r, t) - h_\vartheta(r, t)), H_x(x, y, z, t) = \frac{xz}{r^2} (h_r(r, t) - h_\vartheta(r, t)), \end{aligned} \quad (32)$$

Since the problem is axially symmetric, any plane coming across axis z can be taken to calculate the magnetic field lines. For example, for plane $y = 0$, the differential equation describing the slope of the magnetic field lines looks like

$$\frac{dx}{dz} = \frac{xz(h_r - h_\vartheta)}{r^2 h_\vartheta + z^2(h_r - h_\vartheta)}.$$

The magnetic field lines for the reference solution at time $t = 0.01$ are shown in **Figure 9**.

Results of Eq. (32) are the formulas for the radial and angular components of the field depending on a single space coordinate:

$$\begin{aligned} h_r(r, t) &= H_z(x, y, z, t) + \frac{xH_x(x, y, z, t) + yH_y(x, y, z, t)}{z}, \\ h_\vartheta(r, t) &= H_z(x, y, z, t) - \frac{z(xH_x(x, y, z, t) + yH_y(x, y, z, t))}{r^2 - z^2}, \\ h_\varphi(r, t) &= \frac{xH_y(x, y, z, t) - yH_x(x, y, z, t)}{r} = 0. \end{aligned}$$

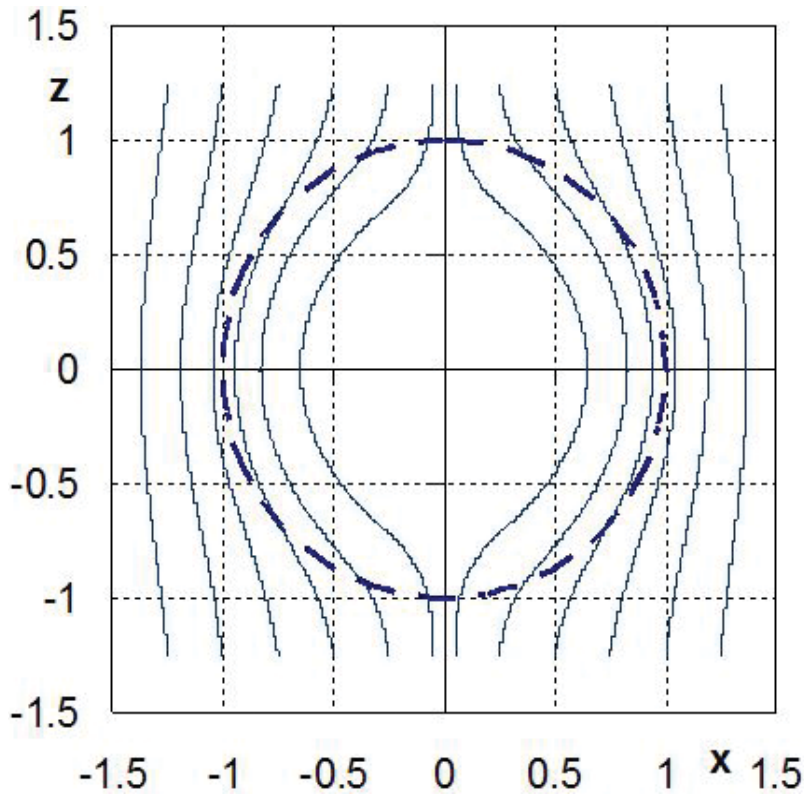


Figure 9. The magnetic field lines in the exact solution with parameters $\nu_1 = 1$ and $\nu_2 \rightarrow \infty$. A dashed line shows the plasma cloud position.

Figures 10 and 11 show profiles of field components for different magnetic viscosity values in external domain ($r > r_0$). Note that there is a small difference between the profiles obtained with $\nu_2 = 50$ and $\nu_2 \rightarrow \infty$ corresponding to the simulation of the plasma cloud interaction with the magnetic field of vacuum $\nu_2 \rightarrow \infty$. For the problem with $\nu_2 \rightarrow \infty$, profiles of the non-dimensionalized field components for the initial phase of diffusion, $t < r_0^3/\nu_1$, are given (see **Figure 12**).

Simulation setup: A computational domain ($|x| \leq 0.5L, |y| \leq 0.5L, |z| \leq 0.5L$) is a cube with edges $L = 10$. Boundary conditions corresponding to the initial undisturbed state are imposed on its lateral faces for the components of field $\mathbf{H}(r, t)|_{r \in \Gamma} = \mathbf{H}(r, t = 0)$. The initial data can be set either for the magnetic field components Eq. (33) or the components of vector potential Eq. (34).

$$\begin{aligned} H_z(x, y, z, t = 0) &= 2f(r, t = 0) + z^2(1 + rf'(r, t = 0))/r^2, \\ H_y(x, y, z, t = 0) &= xyf'(r, t = 0)/r, \quad H_x(x, y, z, t = 0) = xzf'(r, t = 0)/r, \end{aligned} \tag{33}$$

$$\begin{aligned} \Psi_z(x, y, z, t = 0) &= 0, \quad \Psi_y(x, y, z, t = 0) = -zf(r, t = 0), \\ \Psi_x(x, y, z, t = 0) &= yf(r, t = 0). \end{aligned} \tag{34}$$

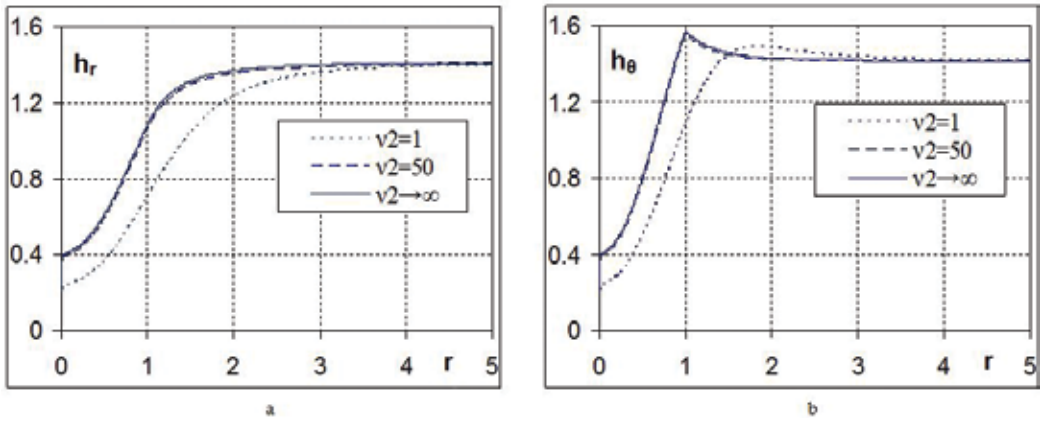


Figure 10. Profiles of the non-dimensionalized field components at time $t = 0.1$: (a) radial h_r and (b) angular h_θ .

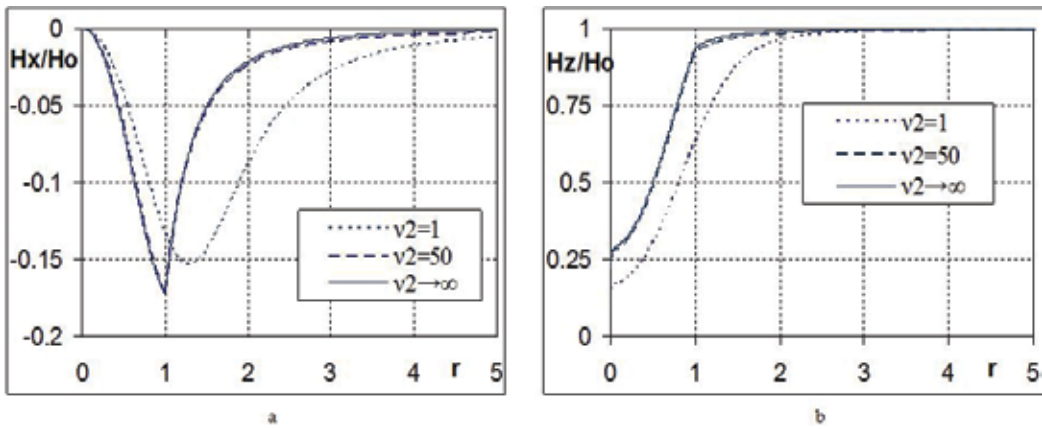


Figure 11. Profiles of the field components along line $x = z$ and $y = 0$ at time $t = 0.1$: (a) H_x and (b) H_z .

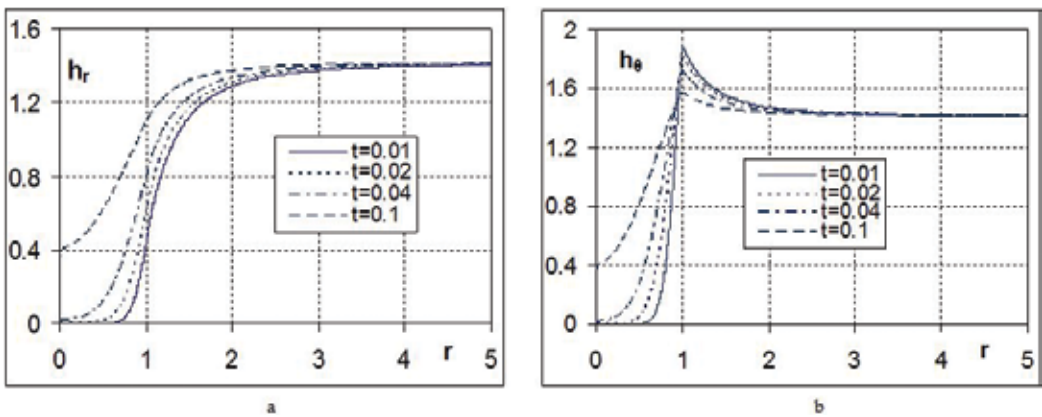


Figure 12. Evolution of non-dimensionalized field components in the problem with parameters $v_1 = 1$ and $v_2 \rightarrow \infty$: (a) radial h_r and (b) angular h_θ .

The EGIDA code uses a scheme preserving the field divergence at one step because difference operators DIV and ROT (div and rot) [18] determined at nodes and in cells of a grid, respectively, satisfy the vector analysis identities: $DIVrot = 0$ and $ROTdiv = 0$.

It has been found that for the first set of initial data the divergence norm depends on errors induced by the initial distribution of the \mathbf{H} field components in the vicinity of sphere $r = r_0$. Though the difference scheme does not change the magnetic field divergence, initial errors lead to a significantly distorted numerical solution (see **Figure 13**).

In the second case, the magnetic field components are determined using the operator numerically differentiating the vector potential, and, hence, the magnetic field divergence norm equals zero at initial time and at all later times. For this case, a good agreement between the calculated results and the exact solution has been achieved even on the coarsest grid (see **Figure 14**).

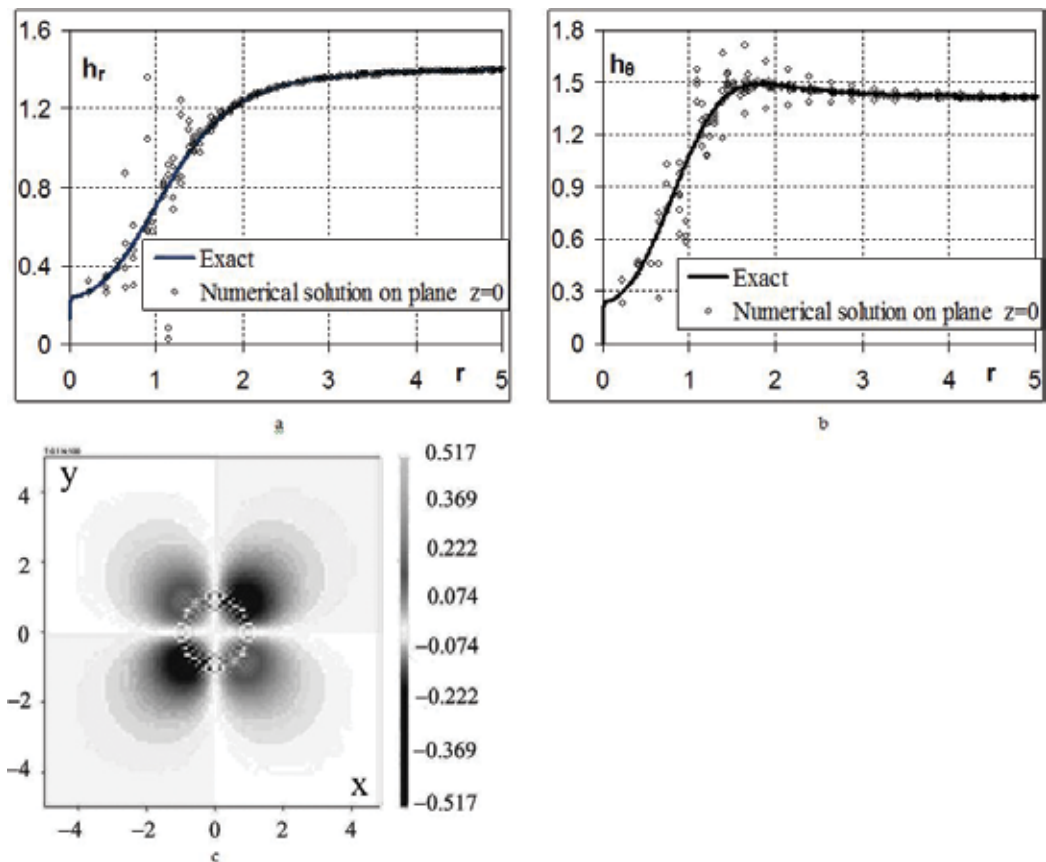


Figure 13. Calculation of the magnetic field diffusion into a spherical plasma cloud for the first set of initial data (33). Profiles of the field components at time $t = 0.1$: (a) radial h_r , and (b) angular h_θ . Distribution of transverse field component H_y in section $z = 0$ (c).

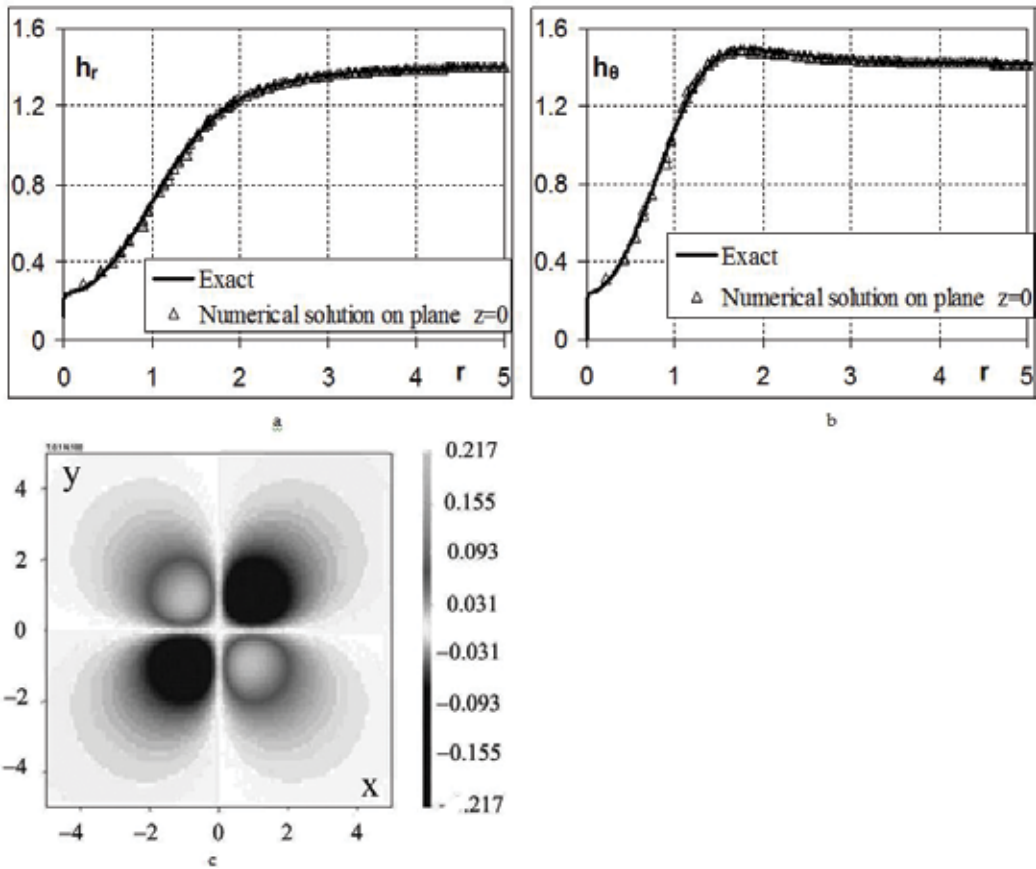


Figure 14. Calculation of the magnetic field diffusion into a spherical plasma cloud for the second set of initial data (34). Profiles of the field components at time $t = 0.1$: (a) radial h_r , and (b) angular h_θ . Distribution of transverse field component H_y in section $z = 0$ (c).

5.2. Diffusion of magnetic field with a low Hall effect

Assume that the Hall effect contribution is small, $bH_0/\nu \ll 1$, but finite. For this reason, it is required to take into account the Hall term in the diffusion equation:

$$\frac{\partial \Psi}{\partial t} = -\nu \cdot \text{rot rot} \Psi - b[\text{rot} \Psi \times \text{rot rot} \Psi].$$

The Hall effect leads to the occurrence of the azimuthal component, H_φ , of magnetic field in plasma [16]:

$$\begin{aligned} H_r(r, \vartheta, t) &= h_r(r, t) \cos \vartheta, H_\vartheta(r, \vartheta, t) = -h_\vartheta(r, t) \sin \vartheta, H_\varphi(r, \vartheta, t) \\ &= h_\varphi \sin \vartheta, h_r(r, t) = 2f(r, t), h_\vartheta(r, t) = \partial r^2 f / r \partial r, h_\varphi(r, t) = \Psi(r, t) / r. \end{aligned}$$

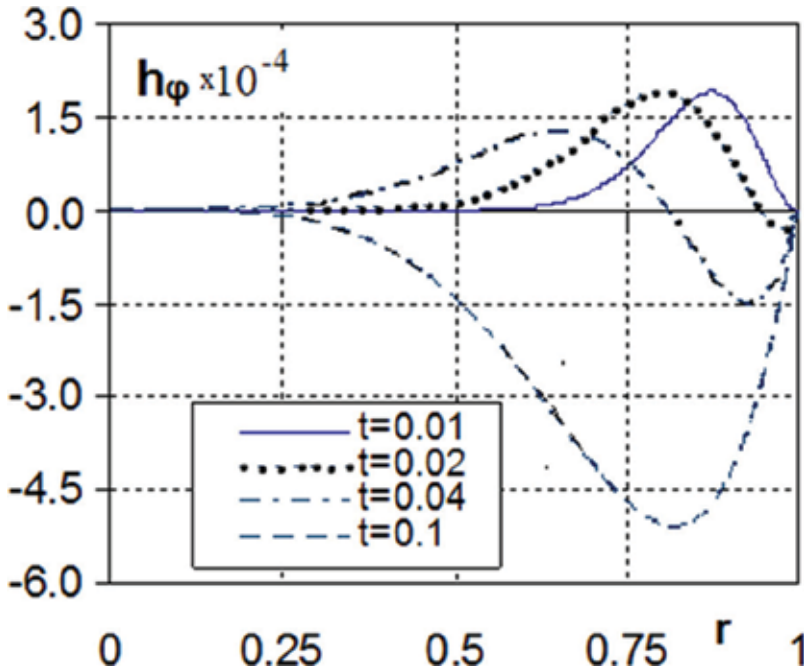


Figure 15. Profiles of the azimuthal component of non-dimensionalized field in the problem with parameters $\nu_1 = 1, \nu_2 \rightarrow \infty, b = 0.01$.

In view of the smallness of parameter bH_0/ν_1 , we have Eqs. (27), (28) to calculate function $f(r, t)$, while the calculation of small additive Ψ is reduced to solving the following boundary value problem:

$$\frac{d\Psi}{dt} = \frac{\partial}{\partial r} \nu \frac{\partial \Psi}{\partial r} - \frac{6\nu\Psi}{r^2} - 2fr^2 \frac{\partial}{\partial r} \left(\frac{b}{r^4} \frac{\partial}{\partial r} r^4 \frac{\partial f}{\partial r} \right), \Psi(r, t = 0) = 0, \Psi(r = 0, t) = 0, \Psi(r \rightarrow \infty, t) = 0.$$

Figure 15 shows profiles of the non-dimensionalized azimuthal field component h_ϕ at early times, $t < r_0^2/\nu_1$. With small values of parameter bH_0/ν_1 , the rest two components $-h_r, h_\theta$ remain unchanged, and they are shown in **Figure 12**. Note that, as it has been shown in [16], if the motion of plasma is accounted, the Hall effect may lead to the occurrence of azimuthal velocity, i.e., to the plasma cloud rotation.

The changeover to Cartesian coordinates is performed using formulas

$$\begin{aligned} H_z(x, y, z, t) &= h_\theta(r, t) + z^2(h_r(r, t) - h_\theta(r, t))/r^2, \\ H_y(x, y, z, t) &= zy(h_r(r, t) - h_\theta(r, t))/r^2 + xh_\phi(r, t)/r, \\ H_x(x, y, z, t) &= xz(h_r(r, t) - h_\theta(r, t))/r^2 - yh_\phi(r, t)/r. \end{aligned}$$

Simulation setup: The initial data is given in the previous section. This problem requires accounting the Hall effect. A local exchange parameter is $b = 0.01$. **Figure 15** shows profiles of the azimuthal component of non-dimensionalized field.

6. Conclusion

An important property of difference schemes in multidimensional flow simulations is that they keep the magnetic field divergence-free in difference solutions. An adverse aspect of this defect is the unphysical transport of matter orthogonal to the field \mathbf{H} [2].

Note that the zero-divergence requirements to difference schemes get more stringent as applied to the solution of diffusion problems. A violation of this requirement results in the accumulation of errors and loss of solution structure, especially in problems with high conductivity gradients.

Acknowledgements

Some EGIDA calculations have been performed under contract no. 1239349 between Sandia National Laboratories and RFNC-VNIIEF. The authors thank sincerely M. Pokoleva for her assistance in simulations and J. Kamm and A. Robinson for their assistance in formulating some benchmarks and references.

Author details

Sofronov Vasily, Zhmailo Vadim and Yanilkin Yury*

*Address all correspondence to: n.yanilkina@mail.ru

VNIIEF, Sarov, Russia

References

- [1] Kulikovskiy AG, Pogorelov NV, Semenov AY. Mathematical Issues of Numerical Solution of Hyperbolic Equations. Moscow: Fizmatlit; 2001 (In Russian)
- [2] Brushlinsky KV. Mathematical and Computational Problems of Magnetic Gas Dynamics. Moscow: Binom; 2009 (in Russian)
- [3] Balsara DS, Spicer DS. A staggered mesh algorithm using high Godunov fluxes to ensure solenoidal magnetic fields in magneto hydrodynamics simulations. *Journal of Computational Physics*. 1999;**149**:270-292
- [4] Toth G. The $\nabla B=0$ constraint in shock-capturing magneto hydrodynamics codes. *Journal of Computational Physics*. 2000;**161**:605-652
- [5] Gardiner TA, Stone JM. An unsplit Godunov method for ideal MHD via constrained transport. *Journal of Computational Physics*. 2005;**205**:509-539

- [6] Dai W, Woodward PR. An approximate Riemann solver for ideal magneto hydrodynamics. *Journal of Computational Physics*. 1994;**111**:354-372
- [7] Brio M, Wu CC. An upwind differencing scheme for the equations of ideal magnetohydrodynamics. *Journal of Computational Physics*. 1988;**75**:400-422
- [8] Takahashi K, Yamada S. Regular and non-regular solutions of the Riemann problem in ideal magnetohydrodynamics. *Journal of Plasma Physics*. 2013;**75**(3):335-356
- [9] Niederhaus JHI. Code verification for ALEGRA using ideal MHD Riemann problems. Sandia National Laboratories preprint SAND2012-1362P. 2008:1-44
- [10] Orszag A, Tang C. Small-scale structure of two-dimensional magneto hydrodynamics turbulence. *Journal of Fluid Mechanics*. 1979;**90**:129-143
- [11] Zhmailo VA, Sofronov VN, Yanilkin YV. Magneto hydrodynamic test problem. *Voprosy Atomnoi Nauki i Tekhniki. Ser. Teoreticheskaya i Prikladnaya Fizika*. 2017;**2**:55-81 (in Russian)
- [12] Yanilkin YV, Belyaev SP, Bondarenko Yu A et al. EGAK and TREK Eulerian codes for multidimensional multimaterial flow simulations. *Transactions of RFNC-VNIIEF. Research publication. Sarov: RFNC-VNIIEF*; 2008. **12**:54-65 (in Russian)
- [13] Eguzhova MY, Zhmailo VA, Sofronov VN, Chernysheva ON, Yanilkin YV, Glazyrin SI. Implementation, analysis and testing of three-dimensional computational methods for MHD simulations of compressible multimaterial flow in Eulerian variables. In: *Proceedings of the 10th Seminar on New Models and Hydro codes for Shock Wave Processes on Condensed Matter*; July 27-August 1, 2014; Pardubice. Czech Republic. pp. 165-176
- [14] Garanin SF, Ivanova GG, Karmishin DV, Sofronov VN. Diffusion of a mega gauss field into a metal. *Journal of Applied Mechanics and Technical Physics*. 2005;**46**(N2):153-159
- [15] Garanin SF. Diffusion of a strong magnetic field into a dense plasma. *Journal of Applied Mechanics and Technical Physics*. 1985;**26**(N3):308-312
- [16] Zhmailo VA, Kokoulin ME. The hall effect in the problem of plasma cloud spread in the magnetic field. *Voprosy Atomnoi Nauki i Tekhniki. Ser. Teoreticheskaya i Prikladnaya Fizika*. 2004;**1-2**:3-12 (in Russian)
- [17] Sedov LI. *Conformity and Dimensionality Methods in Mechanics*. Nauka; 1977 (In Russian)
- [18] Samarsky AA, Tishkin VF, Favorsky AP, Shashkov MY. Operator difference schemes. *Differential Equations*. 1981;**XVII**:N7 (in Russian)

A Numerical Simulation of the Shallow Water Flow on a Complex Topography

Alexander Khoperskov and Sergey Khrapov

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.71026>

Abstract

In current chapter, we have thoroughly described a numerical integration scheme of nonstationary 2D equations of shallow water. The scheme combines the smoothed particle hydrodynamics (SPH) and the total variation diminishing (TVD) methods, which are sequentially used at various steps of the combined SPH-TVD algorithm. The method is conservative and well balanced. It provides stable through calculations in presence of nonstationary “water-dry bottom” boundaries on complex irregular bottom topography including the transition of such a boundary between wet and dry bottom through the computational boundary. Multifarious tests demonstrate the effectiveness of the combined SPH-TVD scheme application for a solution of diverse problems of the engineering hydrology.

Keywords: computational fluid dynamics, numerical methods, hydrology, shallow water, digital elevation model, numerical experiment, GPU

1. Introduction

The usage of the shallow water models (SWMs) allows to solve a wide range of engineering tasks related to the dynamics of surface waters (a seasonal flooding [1], a drain and rain flows [2]), the emergence and expansion of the marine nonlinear waves [3] (problems of tsunami impact on a shore, nonlinear waves formation due to earthquakes, and meteorological waves generation by an open ocean resonance [4]), and the inundation in a coastal area by storm surge [5]. The SWM modifications are effective for studying various geophysical problems such as the dynamics of the pyroclastic flows [6] and the riverbed processes including the sediment dynamics and the diffusion of pollutant particles in reservoirs. The multilayer models utilization significantly expands the opportunities of the shallow water approach [7]. The tasks associated with flooding research of river valleys or interfluves [1] are stood out among the hydrological problems. In the framework of the SWM, the important results have

also been obtained for the atmospheric phenomena, meteorological forecasts, and global climate models [8, 9].

It should be noted that the SWM is also actively applied and developed for theoretical research of various cosmic gas flows. The shallow water approximation allows describing a whole series of astrophysical objects: the protoplanetary and circumstellar disks [10, 11], the accretion disks around the compact relativistic objects [12], the cyclonic movements in the giant planets' atmospheres [13], and the spiral galaxies gas disk components [14]. The gravitational fields play the topography role in such astrophysical problems.

A lot of numerical methods have been proposed for the shallow water dynamics modeling employed for diverse tasks and conditions. Despite the fact that conservative methods of finite volume poorly describe stationary states, they allow correctly calculating shock waves and contact discontinuities. The latter problem could be overcome by the so-called well-balanced (WB) circuits [15–17].

Our main aim is through calculation for a flow with various Froude number within $0 \leq Fr < 100$ ($Fr = u/c_g$, where $c_g = \sqrt{gH}$ is an analog of the sound speed, H is the depth, u is the flow velocity, g is the specific gravitational force) in order to simulate subcritical ($Fr < 1$), transcritical ($Fr \sim 1$), and supercritical ($Fr > 1$) flows. Highly heterogeneous terrain topography including vertical discontinuities and small-scale inhomogeneities at the computational domain boundary makes the calculations noticeably complicated. The latter leads to special quality requirements in numerical algorithms. Hence, the modern numerical schemes should simulate fluid movements along the dry bottom and correctly describe the interfaces between wet and dry bottom [1, 18, 19].

Among the numerous numerical methods solving shallow water equations (SWEs), the following methods should be mentioned: the discontinuous Galerkin method based on triangulation [8], the weighted surface-depth gradient method for the MUSCL scheme [18], and the modified finite difference method [20]. The so-called constrained interpolation profile/multimoment finite-volume method utilizing the shallow water approximation is developed to simulate geophysical currents on a rotating planet in spherical coordinate system ([9] and see the references there). The particle-mesh method demonstrates good opportunities for calculation of rotating shallow water [21]. As a rule, numerical schemes of the second-order accuracy give satisfactory results and allow to correctly solve a wide range of tasks for diverse applications [17]. Special attention should be focused on the numerical way of a source term setting, since in the case of discontinuous topography, it may induce an error at the shock wave front [22].

We consider the original CSPH-TVD (combined smoothed particle hydrodynamics—total variation diminishing) algorithm of numerical integration of the Saint-Venant equations. It accounts for the new modifications improving the computational properties of the scheme. A detailed description of the numerical scheme is the main aim of the chapter.

2. Mathematical models

In this section, the mathematical models for terrestrial hydrology, which accounts for the maximum number of physical and meteorological factors and can be described by the shallow

water equations, are considered. The SWE or the Saint-Venant equations are fairly simple model describing the free surface dynamics of incompressible fluid (see the details in the review of [23]). The model assumes that vertical equilibrium in a medium exists at every moment of time. The conservation laws of mass and momentum in the integral form for a thin layer of moving substance with additional sources and forces are:

$$\frac{d}{dt} \iint_{S(t)} H(\mathbf{r}, t) dS = Q(\mathbf{r}, t), \quad \frac{d}{dt} \iint_{S(t)} H\mathbf{v} dS = -g \iint_{S(t)} H\nabla\eta dS + \iint_{S(t)} H\mathbf{f} dS, \quad (1)$$

where Q is the sources function, $\nabla = \mathbf{e}_x\partial/\partial x + \mathbf{e}_y\partial/\partial y$ is the nabla operator, $S(t)$ is the cross-sectional area of the “liquid particle”, $\mathbf{r} = x\mathbf{e}_x + y\mathbf{e}_y$ is the radius vector, $\mathbf{v} = u\mathbf{e}_x + v\mathbf{e}_y$ is the velocity vector, $\eta(\mathbf{r}, t) = H(\mathbf{r}, t) + b(\mathbf{r})$ is the free surface elevation, $b(\mathbf{r})$ is the bottom profile, $\mathbf{f}(\mathbf{r}, t) = \mathbf{f}_b + \mathbf{f}_{fr} + \mathbf{f}_{Cor} + \mathbf{f}_w + \mathbf{f}_s$ is the sum of the external forces, accounting for the bottom friction \mathbf{f}_b , the internal friction (viscosity) \mathbf{f}_{fr} , the Coriolis force \mathbf{f}_{Cor} , the effect of atmospheric wind \mathbf{f}_w , and the force \mathbf{f}_s determined by the liquid momentum due to the action of the sources Q . The surface density of sources is $\sigma(\mathbf{r}, t) = \sigma^{(+)} + \sigma^{(-)} = dQ/dS$, where $\sigma^{(+)}$ and $\sigma^{(-)}$ are the liquid sources (rain, melting snow, flows through the hydro-constructions, groundwater, etc.) and sinks (infiltration, evaporation, etc.), respectively.

The digital elevation model (DEM) provides the quality of numerical simulations of real hydrological objects to a significant extent. The DEM is determined by the height matrix $b_{ij} = b(x_i, y_j)$ on numerical grid $\{x_i, y_j\}$ ($i = 1, \dots, N_x, j = 1, \dots, N_y$). The DEM elaboration utilizes diverse geoinformation methods for the processing of spatial data obtained from various sources. A matrix of heights has been built in several stages. At the beginning stage, the remote sensing data are accounted by the function $b(x_i, y_j)$. The river sailing directions and the actual water depth measurements allow to construct the DEM for large river beds (for example, for the Volga River and the Akhtuba River). To improve our DEM, the data on various small topography objects such as small waterways, roads, small dams, etc., should be included into the consideration. The numerical simulation results of the shallow water dynamics reveal flooding areas which may be compared with real observational data (both from the remote methods and our own GPS measurements). Such approach qualitatively improves the DEM as a result of the iterative topography refinement.

3. Numerical method CSPH-TVD

In current paragraph, the CSPH-TVD numerical method is thoroughly examined. The method combines classical TVD and SPH methods at various stages of numerical integration of hyperbolic partial differential equations and uses the particular benefits of both. The advantages of graphical processing units (GPUs) with CUDA acceleration for hydrological simulations are also discussed.

The computational domain has been covered by a fixed uniform grid with a spatial step h , while mobile “liquid particles” (hereinafter particles) are placed in the centers of cells. The time layers t_n have a nonuniform step $\tau_n = t_n - t_{n-1}$, where n denotes the index of the time layer. The vector space index $\mathbf{i} = (i, j)$ characterizes the radius vectors of the fixed centers of the cells

$\mathbf{r}_i^0 = (x_i^0, y_j^0)$. The total number of cells and particles is $N_p = N_x \cdot N_y$. At the initial time moment, the particles are placed in the centers of cells. The main characteristics of the particles are the volume $V_i^n = \iint_{S_i(t)} H dS$, momentum $\mathbf{P}_i^n = \iint_{S_i(t)} H v dS$, coordinates of their centers \mathbf{r}_i^n , and linear sizes $\ell_i^n = \sqrt{S_i^n}$.

The CSPH-TVD method main stages are the following:

- I. *The Lagrangian's stage.* It includes an application of the modified CSPH approach [1, 24]. At this stage, the changes of the integral characteristics and particles positions ($V_i^n \rightarrow V_i^{n+1}$, $\mathbf{P}_i^n \rightarrow \mathbf{P}_i^{n+1}$, and $\mathbf{r}_i^n \rightarrow \mathbf{r}_i^{n+1}$) due to the hydrodynamic and external forces acting on them are calculated. **Figure 1** demonstrates the deformations of particles occurred during their motion, while the positions of particles are shifted relatively to the cells centers (it is assumed that $\ell_i^{n+1} \neq h$). The difference between the CSPH-TVD method and the traditional SPH approach [25] consists in the fact that the particles at each time interval $[t_n, t_{n+1}]$ of the numerical time integration at the moment of time t_n are found in the initial state ($\mathbf{r}_i^n \equiv \mathbf{r}_i^0$ and $\ell_i^n = h$).
- II. *The Euler's stage.* At the Euler's stage, a fixed grid is used to calculate the mass and momentum fluxes through the cells boundaries at time moment $t_{n+1/2} = t_n + \tau_n/2$. The corresponding changes of the particles integral characteristics are proportional to the difference between the inflow and outflow fluxes. If the region shaded at **Figure 1** is inside of the cell, then the substance flows into the cell, otherwise it flows out. The modified TVD approach [1, 24] and approximate solution of the Riemann's problem [26] are applied to calculate the flows. At the end of the Euler's stage, the particles return to their initial state.

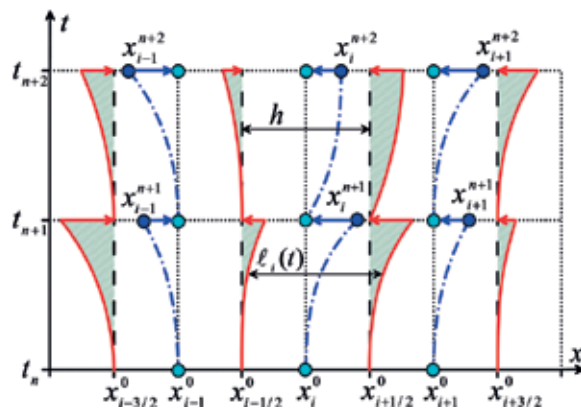


Figure 1. The OX-projection of the main stages scheme of the CSPH-TVD method. The dash-dotted lines $x_i(t)$ correspond to the trajectories of the particles due to their displacement $x_i^n \rightarrow x_i^{n+1}$, the dashed black lines show the boundaries of cells. The particles boundaries deformed during the motion are shown by solid lines. The shaded regions correspond to a change in the particles integral characteristics caused by a flow of matter through the cells boundaries.

It should be noted that when the dry bottom regions are considered, the particles with zero depth and momentum are placed in the corresponding cells. For such particles, the Lagrangian's stage is skipped. At the Euler's stage, the mass and momentum fluxes flowing into the corresponding cell with zero depth are calculated in case of the water presence in the neighboring cells. The changes in the particles integral characteristics are determined too. Such a method permits to carry out an effective through calculation in the presence of nonstationary "water-dry bottom" boundaries in the computational domain.

3.1. The Lagrangian's stage

Let us define the following quantities

$$\varphi_i(t) = \frac{1}{h^2} \iint_{S_i(t)} \varphi(t, x, y) dx dy \quad (2)$$

being the analogous of the function mean values $\varphi = \{H, Hu, Hv\}$ in cells in the finite volume approximation on a fixed grid. Taking into the account Eq. (2), the integral characteristics of the particles can be written in the form:

$$V_i(t) = H_i(t)h^2, \mathbf{P}_i(t) = (Hv)_i(t)h^2 \quad (3)$$

After substituting relations (2) and (3) into the system of Eq. (1), we obtain

$$\frac{d\mathbf{U}_i}{dt} = \mathbf{\Phi}_i \quad (4)$$

where $\mathbf{U}_i = \begin{pmatrix} H_i \\ (Hv)_i \end{pmatrix}$, $\mathbf{\Phi}_i = \begin{pmatrix} \sigma_i \\ -gH_i \nabla_i \eta + H_i \mathbf{f}_i \end{pmatrix}$, $\nabla_i \eta = \nabla H|_{\mathbf{r}=\mathbf{r}_i} + \nabla b|_{\mathbf{r}=\mathbf{r}_i^0}$, $\sigma_i = \sigma(\mathbf{r}_i)$, $\mathbf{f}_i = \mathbf{f}(\mathbf{r}_i)$. The law of motion of a particle is

$$\frac{d\mathbf{r}_i}{dt} = \mathbf{v}_i \quad (5)$$

where $\mathbf{v}_i = (Hv)_i/H_i$ is the velocity of the i -th particle.

For the numerical integration of the system of Eq. (4), an approximation for the spatial derivatives is required. In accordance with the SPH-approach [25], any characteristic of the medium φ and its derivatives $\nabla\varphi$ are replaced by their smoothed values in the flow area Ω :

$$\hat{\varphi}(t, \mathbf{r}) = \sum_{\mathbf{k}=(1,1)}^{(N_x, N_y)} \psi(t, \mathbf{r}_k) W(|\mathbf{r} - \mathbf{r}_k|, h), \nabla \hat{\varphi}(t, \mathbf{r}) = \sum_{\mathbf{k}=(1,1)}^{(N_x, N_y)} \psi(t, \mathbf{r}_k) \nabla W(|\mathbf{r} - \mathbf{r}_k|, h) \quad (6)$$

where W is the smoothing kernel function, h is the smoothing length, \mathbf{k} is the vector index similar to the spatial index \mathbf{i} . Accounting for Eq. (2), the quantity ψ can be determined as

$$\psi(t, \mathbf{r}_k) = \iint_{S_k(t)} \varphi(t, \mathbf{r}) dS = \varphi_k(t) h^2 \tag{7}$$

Our modification of the SPH-method consists in the generalization of Eq. (7). In case of the traditional SPH-approach for the SWE [27] expression $\psi(t, \mathbf{r}_k) = V_k(t) \widehat{\varphi}_k(t) / \widehat{H}_k(t)$ is used instead of (7). Substituting the SPH-approximation (6) and (7) in the right-hand side of the expressions for system (4), we obtain

$$\Phi_i \approx \left(-g H_i \sum_k \left(H_k \nabla_i \overline{W}_{ik} + b_k \nabla_i \overline{W}_{ik}^0 \right) + H_i \mathbf{f}_i \right) \tag{8}$$

where $\overline{W}_{ik} = h^2 W(|\mathbf{r}_i - \mathbf{r}_k|, h)$, $\overline{W}_{ik}^0 = h^2 W(|\mathbf{r}_i^0 - \mathbf{r}_k^0|, h)$ are the functions computed in the cells centers and afterward used to approximate the gradient of the fixed bottom relief $b(\mathbf{r})$. The approximation (8) ensures the conservatism of the CSPH-TVD method at the Lagrangian's stage and the fulfillment of the WB-condition on the inhomogeneous topography.

There are three conditions superimposed on the smoothing kernels W :

1. the kernel is finiteness;
2. the normalization condition is fulfilled $\iint_{\Omega} W(|\mathbf{r} - \mathbf{r}'|, h) d\mathbf{r}' = 1$;
3. the zero spatial step limit $\lim_{h \rightarrow 0} W(|\mathbf{r} - \mathbf{r}'|, h) = \delta(|\mathbf{r} - \mathbf{r}'|, h)$ is the Dirac delta-function.

For the smoothing kernel W , spline functions of different orders or Gaussian distribution are used in the SPH hydrodynamic simulations [25, 28, 29]. The Monaghan's cubic spline

$$W(q) = A_w \begin{cases} 1 - 1.5q^2 + 0.75q^3, & 0 \leq q \leq 1 \\ 0.25(2 - q)^3, & 1 \leq q \leq 2 \\ 0, & q > 2 \end{cases}, W'(q) = -\frac{A_w}{h} \begin{cases} 3q - 2.25q^2, & 0 \leq q \leq 1 \\ 0.75(2 - q)^2, & 1 \leq q \leq 2 \\ 0, & q > 2 \end{cases} \tag{9}$$

is the most commonly applied approximation, where $q = |\mathbf{r}_i - \mathbf{r}_k| / h$ is the relative distance between the particles, $A_w = \left(\frac{2}{3h}, \frac{10}{7\pi h^2}, \frac{1}{\pi h^3} \right)$ is the normalization constant for the 1D, 2D and 3D cases, respectively. In the CSPH-TVD method, the smoothing kernel W approximates only the gradients of physical quantities in Eq. (8). Therefore, the value of the normalization constant may be corrected to increase the accuracy of the numerical solution. For example, it was done by the authors of [1] when for the two-dimensional case, the normalization constant was chosen to be $A_w \approx 0.987 \left(\frac{10}{7\pi h^2} \right)$. The smoothing kernel (9) can lead to unphysical clustering of particles in the simulation of hydrodynamic flows by the SPH method; therefore, an approximation for the smoothing kernel of the pressure gradient is applied in form [29]:

$$W(q) = B_w \begin{cases} (2-q)^3, & 0 \leq q \leq 2 \\ 0, & q > 2 \end{cases}, W'(q) = -\frac{B_w}{h} \begin{cases} 3(2-q)^2, & 0 \leq q \leq 2 \\ 0, & q > 2 \end{cases} \quad (10)$$

where $B_w = \left(\frac{1}{8h}, \frac{5}{16\pi h^2}, \frac{15}{64\pi h^3}\right)$ is the normalization constant. **Figure 2** shows the comparison between two smoothing kernels (9) and (10). As seen from **Figure 2b** for $q < 0.7$, the characteristic feature of the smoothing kernel (9) is the impetuous decrease of the derivative value to zero as the particles approach each other. At small q values, this feature leads to a significant weakening of the hydrodynamic repulsion force between the particles and appearance of nonphysical clusters of particles.

For the numerical integration of the set of differential Eqs. (4) and (5), the method of a predictor-corrector type of second-order accuracy (so-called leapfrog method) is employed. The main steps of the leapfrog method for the Lagrange stage of CSPH-TVD method are:

- I. At the predictor step, the water depth H_i and velocity \mathbf{v}_i are calculated at time moment $t_{n+1} = t_n + \tau_n$:

$$\mathbf{U}_i^*(t_{n+1}) = \mathbf{U}_i(t_n) + \tau_n \Phi_i(\mathbf{r}_i(t_n), \mathbf{U}_i(t_n)) \quad (11)$$

- II. The Particles' spatial positions \mathbf{r}_i are determined at time t_{n+1} :

$$\mathbf{r}_i(t_{n+1}) = \mathbf{r}_i(t_n) + \frac{\tau_n}{2} (\mathbf{v}_i(t_n) + \mathbf{v}_i^*(t_{n+1})) \quad (12)$$

- III. During the corrector step, the water depth H_i and velocity \mathbf{v}_i values are recalculated (t_{n+1}):

$$\tilde{\mathbf{U}}_i(t_{n+1}) = \frac{\mathbf{U}_i(t_n) + \mathbf{U}_i^*(t_{n+1})}{2} + \frac{\tau_n}{2} \Phi_i(\mathbf{r}_i(t_{n+1}), \mathbf{U}_i^*(t_{n+1})) \quad (13)$$

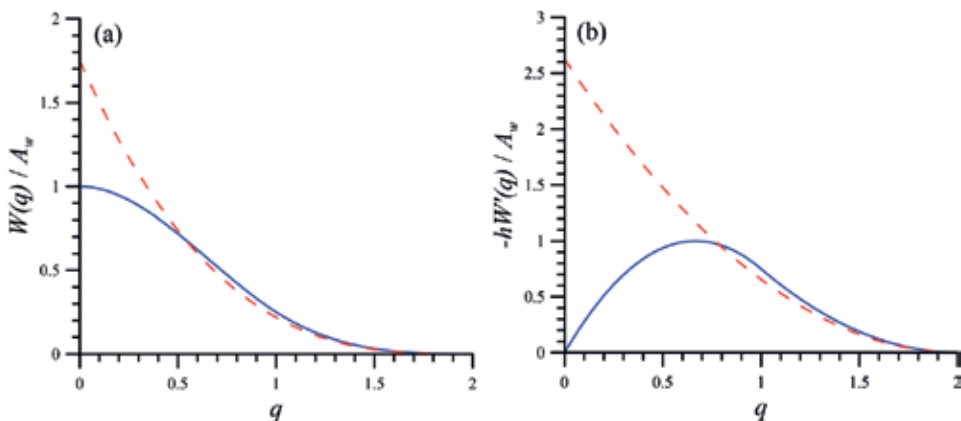


Figure 2. The spatial distributions of smoothing kernels (a) and their derivatives (b). The smoothing kernels (9) and (10) are shown by solid and dashed lines, correspondingly.

To reduce the rounding errors during the numerical implementation of the algorithm in recurrent relations (11)–(13), the value of the relative displacement of particle $\xi_i(t) = \mathbf{r}_i(t) - \mathbf{r}_i^0$ (here \mathbf{r}_i^0 is the initial position) should be used instead of $\mathbf{r}_i(t)$.

3.2. The Euler’s stage

At this stage, the relationship between the values of $\tilde{\mathbf{U}}_i^{n+1}$ and \mathbf{U}_i^{n+1} is determined at time moment t_{n+1} . According to the CSPH-TVD approach, the changes in the particles integral characteristics at their return to the initial state \mathbf{r}_i^0 are computed. In order to do it, the following expression

$$\frac{d}{dt} \iint_{S_i(t)} \mathbf{U} dx dy = \iint_{S_i(t)} \left(\frac{\partial \mathbf{U}}{\partial t} + \frac{\partial \mathbf{F}}{\partial x} + \frac{\partial \mathbf{G}}{\partial y} \right) dx dy, \text{ where } \mathbf{F} = \begin{pmatrix} H v \\ H u^2 \\ H v u \end{pmatrix}, \mathbf{G} = \begin{pmatrix} H v \\ H u v \\ H v^2 \end{pmatrix}$$

should be integrated over time from t_n to t_{n+1} . Taking into the account Eq. (2), we obtain

$$\mathbf{U}_i^{n+1} = \tilde{\mathbf{U}}_i^{n+1} - \frac{\tau_n}{h} \left(\mathbf{F}_{i+1/2,j}^{n+1/2} - \mathbf{F}_{i-1/2,j}^{n+1/2} + \mathbf{G}_{i,j+1/2}^{n+1/2} - \mathbf{G}_{i,j-1/2}^{n+1/2} \right) \tag{14}$$

where $\mathbf{F}_{i\pm 1/2,j}^{n+1/2}$ and $\mathbf{G}_{i\pm 1/2,j}^{n+1/2}$ are the average values (in the interval $[t_n, t_{n+1})$) of mass and momentum fluxes through the cell boundaries in x and y directions, respectively. To determine these fluxes, we use the TVD-approach [30] and approximate methods of the Riemann’s Problem (RP) solving for an arbitrary discontinuity decay (Lax-Friedrichs (LF) method, Harten-Lax-van Leer (HLL) method [31]). The Riemann’s problem is solved separately for each boundary of the Euler’s cells.

The values of the flow parameters to the left \mathbf{U}^L and right \mathbf{U}^R of the considered boundary are the initial conditions for RP. These quantities define magnitude of the jump and can be obtained by a piecewise polynomial reconstruction of the function $\mathbf{U}(t, \mathbf{r})$. We limit the consideration by the piecewise linear reconstruction providing the second order of accuracy by space coordinate for a numerical scheme [32]. The modification of the TVD-approach used in the CSPH-TVD method is concerned with the reconstruction of the value of $\mathbf{U}(t, \mathbf{r})$ relatively to the position of the particles at time $t_{n+1/2}$ at distance $\xi_i(t_{n+1/2})$ from the centers of the cells. Thus, the piecewise linear profile of the function $\mathbf{U}(t, \mathbf{r})$ inside of the i -th cell at time $t_{n+1/2}$ in the projection onto the Ox axis (**Figure 3**) has the following form

$$\mathbf{U}(t_{n+1/2}, x) = \tilde{\mathbf{U}}_{i,j}^{n+1/2} + \left(x - x_i^{n+1/2} \right) \Theta_{i,j}^{n+1/2}, \quad x \in \left[x_{i-1/2}^0, x_{i+1/2}^0 \right] \tag{15}$$

where $\Theta_{i,j}^{n+1/2}$ is the vector of slopes (angular coefficients), $x_i^{n+1/2} = x_i^0 + \xi_i^{n+1/2}$. Taking into the account Eq. (15), the expressions for \mathbf{U}^L and \mathbf{U}^R at the boundary $x_{i+1/2}^0$ can be written as

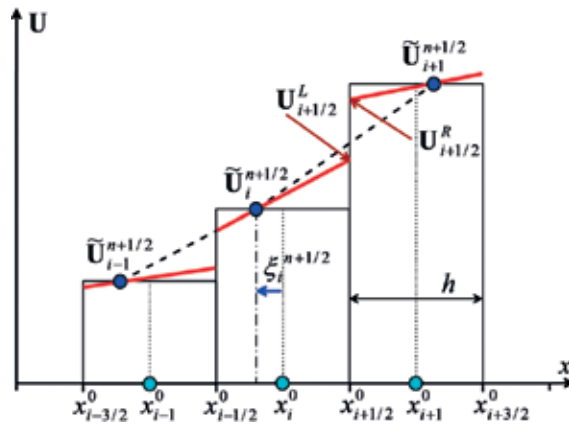


Figure 3. The projection onto the Ox axis of the piecewise linear reconstruction of the function $U(t, \mathbf{r})$ with the respect to the mass centers of the particles for the CSPH-TVD scheme at time $t_{n+1/2}$.

$$U_{i+1/2}^L = \tilde{U}_{i,j}^{n+1/2} + \left(\frac{h}{2} - \xi_i^{n+1/2}\right) \Theta_{i,j}^{n+1/2}, U_{i+1/2}^R = \tilde{U}_{i+1,j}^{n+1/2} - \left(\frac{h}{2} + \xi_{i+1}^{n+1/2}\right) \Theta_{i+1,j}^{n+1/2} \quad (16)$$

The slopes of the piecewise linear distribution (16) should satisfy the TVD-condition [30]. To fulfill this condition, the limiter function

$$\Theta_{i,j}^{n+1/2} = \Lambda \left(\frac{\tilde{U}_{i+1}^{n+1/2} - \tilde{U}_i^{n+1/2}}{h + \xi_{i+1}^{n+1/2} - \xi_i^{n+1/2}}, \frac{\tilde{U}_i^{n+1/2} - \tilde{U}_{i-1}^{n+1/2}}{h + \xi_i^{n+1/2} - \xi_{i-1}^{n+1/2}} \right)$$

has to be used. The analysis shows that the limiters of the minmod [33], van Leer [32], and van Albada et al. [34] satisfy the TVD-condition and suppress the nonphysical oscillations of the numerical solution in the vicinity of discontinuities.

3.3. The stability condition and parallel implementation on GPUs

For the stability of the numerical scheme CSPH-TVD the following conditions have to be fulfilled during the integration time τ_n : (1) at the Lagrangian's stage the shift of particles' center of mass should not exceed $h/2$ relatively to the initial position; (2) at the Euler's stage the perturbations should not propagate over a distance greater than the cell size h . These requirements provide the numerical stability condition for the time step τ_n of the CSPH-TVD algorithm:

$$\tau_n = K \min \left(\frac{h}{2 \max_i |\mathbf{v}_i^n|}, \frac{h}{\max_i (|\mathbf{v}_i^n| + \sqrt{gH_i^n})} \right) \quad (17)$$

where $0 < K < 1$ is the Courant number.

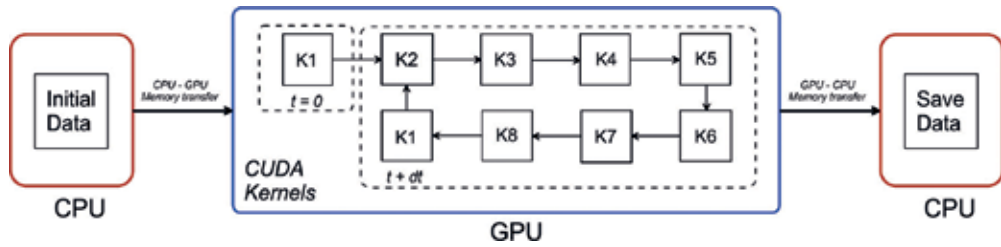


Figure 4. The flow diagram for each of the calculation module: K1 determines the presence of water in the CUDA block; K2 calculates the forces at the time moment t_n at the Lagrangian’s stage; K3 calculates the time step t_{n+1} ; K4 calculates the new positions of the particles and their integral characteristics at time $t_{n+1/2}$; K5 defines the forces on the time layer $t_{n+1/2}$; K6 calculates the positions of the particles and their integrated characteristics for the next time layer t_{n+1} ; K7 calculates the flux of physical quantities through the cells boundaries at time moment $t_{n+1/2}$; and K8 determines the final hydrodynamic parameters at time moment t_{n+1} (see details in Ref. [35]).

The CSPH-TVD method demonstrates a good ability for parallelization on graphics processors [35]. **Figure 4** shows the computer implementation of the CSPH-TVD method on GPUs. It indicates the execution sequence of CUDA kernels at various stages of our algorithm.

4. The set of test problems verifying the numerical models

In the paragraph, the basic set of various tests for one- and two-dimensional problems suitable for revelation of positive and negative sides of the numerical shallow water models is reviewed. According to Toro [26], there are four main test types appropriate to evaluate the numerical solutions (**Figure 5**).

1. The comparison with the exact solutions for several 1D problems.
 - a. Dam-break waves propagating on the wet bed [16, 36] with the initial distribution of the water depth $H(x \leq 0) = H_{01}$ and $H(x > 0) = H_{02}$, ($H_{01} > 0, H_{02} > 0$). This solution is an analog of the pressure jump decay in an ordinary hydrodynamics. It also contains a hydraulic shock (an analog of a shock wave) and a rarefaction wave for which the exact solutions of the RP for a self-similar type $f(x, t) = f(\xi) \equiv f(x/t)$ exist [37]. The bora (hydraulic jump) front for the CSPH-TVD scheme is smeared into three cells and does not contain dispersion oscillations typical for sufficiently good numerical schemes.
 - b. Dam-break waves propagating on the dry bed [16]: $H(x \leq 0) = H_{01}$ and $H(x > 0) = 0$. This is an important test imposing special requirements on a numerical scheme for a

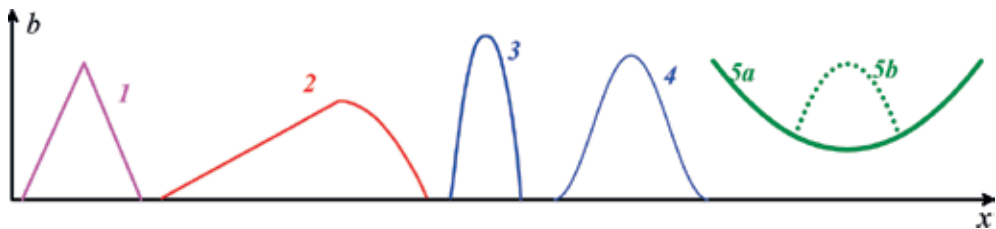


Figure 5. The bottom profiles with known exact solutions of SWE.

correct description of the moving boundary between the water and dry bottom. It also allows comparison of the numerical result with exact solution of the RP [37].

- c. The dam-break problem in a sloping channel is more complicated test comparing the dynamics of numerical wetting and drying fronts with the analytic solutions [18].
 - d. The two shock waves collision is described in the terms of the RP and contains a rarefaction wave [36].
 - e. It is possible to check the coincidence of the numerical and exact solutions for the flow over the bed profiles of different shapes (**Figure 5**). The transcritical mode with and without a shock for the flow over bed profile in the form of a parabolic bump (see line 3, **Figure 5**) [16–18, 38] or triangular obstacle with a break (see line 1) [24] have been studied.
 - f. The more complicated bottom profile (see line 2 in **Figure 5**) helps identifying the best advantages of numerical schemes [22]. We have exact stationary solutions on a stepped bottom for subcritical (Froude number $Fr < 1$) and supercritical regimes ($Fr > 1$) [24, 39].
 - g. To verify the properties of well-balanced numerical model, the oscillating lake with a heterogeneous bottom [15] or the basin with a stationary fluid with an island in the middle of the numerical domain are suitable to modeling [19]. The study of a propagation of small perturbation in stationary water with a nonuniform bottom allows to determine the quality of the well-balanced scheme [16, 19] and guarantees the absence of a numerical storm produced by the scheme [40].
 - h. The same problem of assessing the quality of the WB-scheme is solved by studying the oscillations of the liquid level on the 1D parabolic bottom [24], when the profile of the oscillating free surface of the liquid should remain flat at any time. The two-dimensional analog of such problem is the oscillations in the frictionless paraboloid basin $b(x, y) = b_L(x^2 + y^2)/L^2$, for which the free surface remains flat [18, 41]. The analytical solution for the Thacker test with curved water surface in the paraboloidic bottom basin is described in [41]. The initial circular paraboloidic water volume oscillates periodically changing the shape of the free surface in a complex manner (see **Figure 5**, lines **5a** and **5b** corresponds to the bottom profile and the initial shape of the water volume) [18].
2. The second approach is based on the tests with reliable numerical solutions for 1D equations. A simple and effective two-dimensional test is circular dam-break [16, 36], since for the exact solution can be obtained from the one-dimensional radial equation in the polar coordinate system [18]. It is necessary to distinguish the circular dam break on the wet and dry bed.
 3. The third direction is represented by tests using other numerical schemes.
 4. Eventually, the numerical SWM solutions can be compared with the observables data or 3D numerical solutions.
 - a. Dam-break flow over the initially dry bed with a bottom obstacle is a tough test for a simple SWM due to the formation of negative wave propagation [42].

- b. The experimental results of the dam-break on a dry bed channel with varying width are used to examine numerical models [16]. The laboratory measurements of flow parameters for the dam-break in the nonprismatic converging-diverging channel are described in Ref. [38].
- c. The so-called right-angled open channel junction test gives a complex two-dimensional flow structure which is highly inhomogeneous in the cross section of the main channel [43]. Such a kind of comparison with essentially 2D currents is an important stage in testing of the numerical model [16].
- d. The dam-break flow over a triangular obstacle contains a large number of characteristic features [38], including the downstream over the dry channel bed, the smooth rarefaction wave, the upstream reflected shock, and the downstream shock after passing the triangular obstacle. The formation of a reflected shock wave from the boundary wall passing back through the hump causes complex nonlinear oscillations in such a basin.
- e. An analytical solution for steady flow case with spatially varying width, bed, and friction slope of special kind is known [38, 44]. There is a singular point in the flow in the transition from the inclined profile to the shallow bottom. A fine structure of the solution in the vicinity of this point is a good test for numerical models.
- f. The use of various software, for example, FLOW-3D, FLOW-3D-SWM, MIS2D, is a fast and effective way for the verification of new numerical models of shallow water [45].

Figure 6 shows the results of the numerical simulations of shear flow instability (Kelvin-Helmholtz instability) in shallow water using the CSPH-TVD method. The characteristic stages

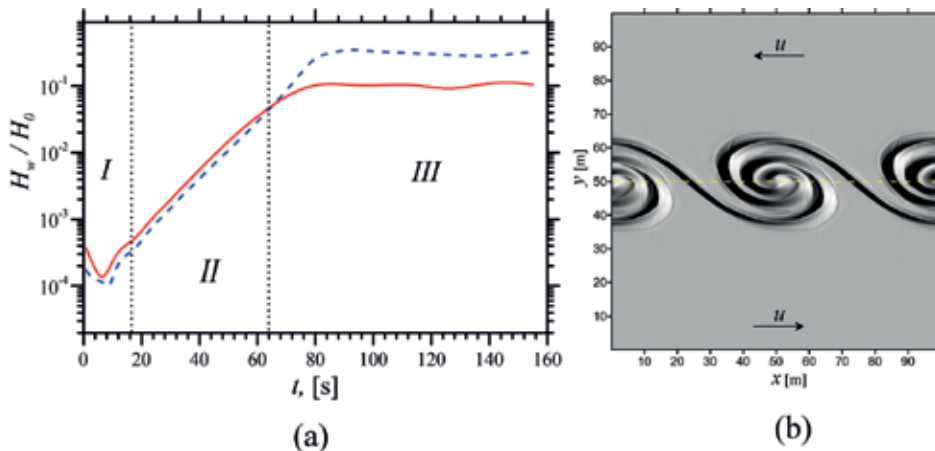


Figure 6. The problem of the tangential velocity discontinuity is represented by: (a) the time dependence of the minimum (shown by solid line) and maximum (shown by dashed line) of the amplitude of the water depth disturbances $H_w(t) = H(t) - H_0$ ($H_0 \equiv H(0)$); (b) the spatial distribution of the vorticity of the velocity vector field $\omega = \nabla \times v/2$ at $t = 110$. The dashed line shown the position of the tangential discontinuity.

of the instability development of the tangential velocity discontinuity in shallow water (see **Figure 6a**) are obvious: (I) the formation of the eigenmode; (II) the linear stage of instability development, when perturbations increase according to the law $\propto \exp(t/T)$ ($T \approx 9.6$ for $Fr = 1$); (III) the nonlinear mode of perturbations evolution, when the growth of the perturbation amplitude ceases and the characteristic vortex structures are formed (so-called cat eyes, **Figure 6b**). Thus, in the CSPH-TVD method, there is a lack of important disadvantage of classical SPH algorithm, which does not allow correctly modeling the tangential discontinuity.

5. Engineering applications

We have developed software to solve some engineering applications utilizing the parallel implementation of the CSPH-TVD method on GPUs.

1. The hydrological regimes of the spring flooding in the Volga-Akhtuba interfluve have been studied thoroughly [1, 35]. The optimal hydrographs of the Volga Hydroelectric Power Station have been constructed. And finally a new approach optimizing the hydrotechnical projects has been developed [46–48].
2. Our numerical experiments have reproduced the dynamics of the catastrophic flooding wave in the city of Krymsk area (Russia, Caucasus Mountains) in July 2012, leading to massive losses of life. A number of features of the hydrological regime during the flash flood of 2012 [49], associated with the landscape and the distribution of rainfall has been revealed.
3. **Figure 7** shows the more detailed results of the numerical simulations of the tsunami formation in the Pacific Ocean during an earthquake off the coast of Japan in 2011. A bilateral hydraulic shock (tsunami) is formed due to the displacement of tectonic plates and then the waves propagate to both sides of the discontinuity. The tsunami reaches the shores of Japan in 20–40 min after the earthquake, while the wave height H_w reaches 10–20 m at the coast line.

6. Conclusion

We have described in detail the numerical method for solving the equations of hydrodynamics in the approximation of the SWM. Our method is a hybrid scheme successfully combining positive properties of Euler's and Lagrange's approaches. The CSPH-TVD method allows to model nonstationary flows on a complex heterogeneous bottom relief, containing kinks and sharp jumps of the bottom profile. The numerical scheme is conservative, well balanced and provides a stable through calculation in the presence of non-stationary "water-dry bottom" boundaries on the irregular bottom relief, including the transition through the computing boundary.

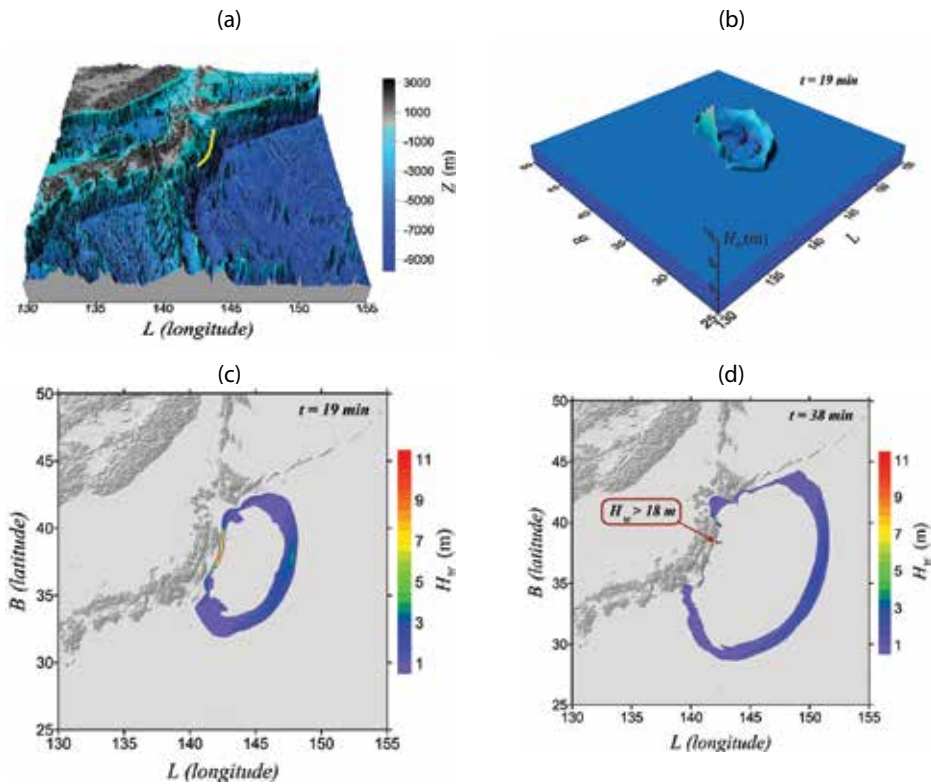


Figure 7. Numerical modeling of the tsunami occurred in 2011 off the coast of Japan: (a) the digital terrain model, where the size of the simulation area is 1500×1500 cells; the yellow line shows the area of tsunami formation; (b) the 3D wave structure 19 min after the earthquake; frames (c) and (d) are the tsunamis at different moments of time.

On the SPH stage, various smoothing cores can be applied, as well as various TVD-delimiters and methods for the RP solution depending on the features of the problem being solved. The scheme has the second order of convergence on smooth solutions and the first order on discontinuities that corresponds to the accuracy of the Godunov-type schemes. In the case of a non-uniform topography the CSPH-TVD method requires less computational resources than the Godunov's type schemes, when the WB-condition is necessary to fulfill for the numerical scheme. Comparing with various SPH-method modifications, the numerical CSPH-TVD scheme has higher accuracy and computational speed for the same number of particles; it is less dissipative and better balanced.

Acknowledgements

The first author has been supported by the Ministry of Education and Science of the Russian Federation (government task No. 2.852.2017/4.6). The second author is thankful to the RFBR (grants 16-07-01037, 15-45-02655, 16-45-340152 and 16-48-340147).

Author details

Alexander Khoperskov* and Sergey Khrapov

*Address all correspondence to: khoperskov@volsu.ru

Volgograd State University, Volgograd, Russia

References

- [1] Khrapov S, Pisarev A, Kobelev I, Zhumaliev A, Agafonnikova E, Losev A, Khoperskov A. The numerical simulation of shallow water: Estimation of the roughness coefficient on the flood stage. *Advances in Mechanical Engineering*. 2013 Id. 787016;5:1-11
- [2] Singh J, Altinakar MS, Ding Y. Numerical modeling of rainfall-generated overland flow using nonlinear shallow-water equations. *Journal of Hydrologic Engineering*. 2015;20 Id. 04014089
- [3] Kowalik Z. *Introduction to Numerical Modeling of Tsunami Waves*. Fairbank: University of Alaska; 2012 196 p
- [4] Sepic J, Vilibic I, Fine I. Northern Adriatic meteorological tsunamis: Assessment of their potential through ocean modeling experiments. *Journal of Geophysical Research: Oceans*. 2015;120:2993-3010
- [5] Jeong W. A study on simulation of flood inundation in a coastal urban area using a two-dimensional well-balanced finite volume model. *Natural Hazards*. 2015;77:337-354
- [6] Mangeney A, Bouchut F, Thomas N, Vilotte JP, Bristeau MO. Numerical modeling of self-channeling granular flows and of their levee-channel deposits. *Journal of Geophysical Research*. 2007 id. F02017;112:1-21
- [7] Chertock A, Kurganov A, Qu Z, Wu T. Three-layer approximation of two-layer shallow water equations. *Mathematical Modelling and Analysis*. 2013;18:675-693
- [8] Lauter M, Giraldo FX, Handorf D, Dethloff K. A discontinuous galerkin method for the shallow water equations in spherical triangular coordinates. *Journal of Computational Physics*. 2008;227:10226-10242
- [9] Li X, Chen D, Peng X, Takahashi K. A multimoment finite-volume shallow-water model on the yin–yang overset spherical grid. *Monthly Weather Review*. 2008;136:3066-3086
- [10] Khoperskov AV, Khrapov SS, Nedugova EA. Dissipative-acoustic instability in accretion disks at a nonlinear stage. *Astronomy Letters*. 2003;29(4):246-257
- [11] Regaly Z, Vorobyov E. The circumstellar disk response to the motion of the host star. *Astronomy and Astrophysics*. 2017;601:A24

- [12] Donmez O. The effects of different perturbations on the dynamics of the accretion disk around the black hole. *International Journal of Modern Physics A*. 2007;**22**:1875-1898
- [13] Nezlin MV, Snezhkin EN. *Rossby Vortices, Spiral Structures, Solitons*. Springer-Verlag; 1993 223 p
- [14] Zasov AV, Saburova AS, Khoperskov AV, Khoperskov SA. Dark matter in galaxies. *Physics-Uspexhi*. 2017;**60**:3-39
- [15] Audusse E, Bouchut F, Bristeau M-O, Perthame LB. A fast and stable well-balanced scheme with hydrostatic reconstruction for shallow water flows. *SIAM Journal on Scientific Computing*. 2004;**25**:2050-2065
- [16] Amiri SM, Talebbeydokhti N, Baghlani A. A two-dimensional well-balanced numerical model for shallow water equations. *Scientia Iranica*. 2013;**20**:97-107
- [17] Audusse E, Bouchut F, Bristeau M-O, Sainte-Marie J. Kinetic entropy inequality and hydrostatic reconstruction scheme for the Saint-Venant system. *Mathematics of Computation*. 2016;**85**:2815-2837
- [18] Aureli F, Maranzoni A, Mignosa P, Ziveri C. A weighted surface-depth gradient method for the numerical integration of the 2D shallow water equations with topography. *Advances in Water Resources*. 2008;**31**:962-974
- [19] Vater S, Beisiegel N, Behrens J. A limiter-based well-balanced discontinuous galerkin method for shallow-water flows with wetting and drying: One-dimensional case. *Advances in Water Resources*. 2015;**85**:1-13
- [20] Zarmehi F, Tavakoli A. A simple scheme to solve saint-venant equations by finite element method. *International Journal of Computational Methods*. 2016 id. 1650001;**13**:1-22
- [21] Frank J, Gottwald G, Reich S. A hamiltonian particle-mesh method for the rotating shallow-water equations. *Lecture Notes in Computational Science and Engineering*. 2003;**26**:131-142
- [22] Navas-Montilla A, Murillo J. Overcoming numerical shockwave anomalies using energy balanced numerical schemes. Application to the Shallow Water Equations with discontinuous topography. *Journal of Computational Physics*. 2017;**340**:575-616
- [23] Zeytounian RK. Nonlinear long waves on water and solutions. *Physics-Uspexhi*. 1995;**38**: 1333-1381
- [24] Khrapov SS, Khoperskov AV, Kuz'min NM, Pisarev AV, Kobelev IA. A numerical scheme for simulating the dynamics of surface water on the basis of the combined SPH-TVD approach. *Numerical Methods and Programming*. 2011;**12**:282-297
- [25] Monaghan J. Smoothed Particle Hydrodynamics. *J. Annual Review of Astronomy and Astrophysics*. 1992;**30**:543-574
- [26] Toro EF. *Riemann Solvers and Numerical Methods for Fluid Dynamics: A Practical Introduction*. 3rd ed. Berlin Heidelberg: Springer-Verlag; 2009 721 p

- [27] Monaghan J. Smoothed particle hydrodynamics. *Journal of Report on Progress in Physics*. 2005;**68**:1703-1759
- [28] Desbrun M, Cani M-P. Smoothed particles: A new paradigm for animating highly deformable bodies. In: *Proceedings of the Eurographics Workshop on Computer Animation and Simulation*. 1996. p. 61-76
- [29] Muller M, Charypar D, Gross M. Particle-based fluid simulation for interactive applications. *Proceedings of 2003 ACM SIGGRAPH Symposium on Computer Animation*. 2003:154-159
- [30] Harten A. High resolution schemes for hyperbolic conservation laws. *Journal of Computational Physics*. 1983;**49**:357-393
- [31] Harten A, Lax P, van B. On upstream differencing and Godunov type methods for hyperbolic conservation laws. *SIAM Review*. 1983;**25**:35-61
- [32] van B. Towards the ultimate conservation difference scheme V, a second order sequel to Godunov's method. *Journal of Computational Physics*. 1979;**32**:110-136
- [33] Roe PL. Some contributions to the modelling of discontinuous flows. *Proceedings of the SIAM/AMS Seminar*. 1983:163-193
- [34] van Albada GD, van B, Roberts WW. A comparative study of computational methods in cosmic gas dynamics. *Astronomy and Astrophysics*. 1982;**108**(1):76-84
- [35] Dyakonova T, Khoperskov A, Khrapov S. Numerical model of shallow water: The use of NVIDIA CUDA graphics processors. *Communications in Computer and Information Science*. 2016;**687**:132-145
- [36] LeVeque RJ. *Finite Volume Methods for Hyperbolic Problems*. UK: Cambridge University Press; 2004 558 p
- [37] Kulikovskii AG, Pogorelov NV, Semenov AY. *Mathematical Aspects of Numerical Solution of Hyperbolic Systems*. Chapman & Hall/CRC; 2000. p. 560
- [38] Alias NA, Liang Q, Kesserwani G. A Godunov-type scheme for modeling 1D channel flow with varying width and topography. *Computers & Fluids*. 2011;**46**:88-93
- [39] Alcrudo F, Benkhaldoun F. Exact solutions to the Riemann problem of the shallow water equations with a bottom step. *Computers & Fluids*. 2001;**30**(6):643-671
- [40] Noelle S, Pankratz N, Puppo G, Natvig JR. Well-balanced finite volume schemes of arbitrary order of accuracy for shallow water flows. *Journal of Computational Physics*. 2005;**213**:474-499
- [41] Thacker WC. Some exact solutions to the nonlinear shallow-water wave equations. *Journal of Fluid Mechanics*. 1981;**107**:499-508
- [42] Ozmen-Cagatay H, Kocaman S. Dam-break flow in the presence of obstacle: Experiment and CFD simulation. *Engineering Applications of Computational Fluid Mechanics*. 2011;**5**: 541-552

- [43] Weber LJ, Schumate ED, Mawer N. Experiments on flow at a 90° open-channel junction. *ASCE Journal of Hydraulic Engineering*. 2001;**127**(5):340-350
- [44] Birman A, Falcovitz J. Application of the GRP scheme to open channel flow equations. *Journal of Computational Physics*. 2007;**222**:131-154
- [45] Wang, T. and Chu, V., Manning friction in steep open-channel flow. *Seventh International Conference on Computational Fluid Dynamics, ICCFD7–3303* (2012), 14 p
- [46] Vasilchenko A, Voronin A, Svetlov A, Antonyan N. Assessment of the impact of riverbeds depth in the northern part of the Volga-Akhtuba floodplain on the dynamics of its flooding. *International Journal of Pure and Applied Mathematics*; **110**(1):183-192
- [47] Voronin AA, Vasilchenko AA, Pisarev AV, Khrapov SS, Radchenko YE. Designing mechanisms of the hydrological regime management of the volga-akhtuba floodplain based on geoinformation and hydrodynamic modeling. *Science Journal of VolSU. Mathematics. Physics*. 2016;**1**:24-37
- [48] Voronin A, Vasilchenko A, Pisareva M, Pisarev A, Khoperskov A, Khrapov S, Podshipkova J. Designing a system for ecological–economical management of the Volga–Akhtuba floodplain on basis of hydrodynamic and geoinformational simulation. *Upravlenie bol'simi sistemami*. 2015;**55**:79-102
- [49] Agafonnikova EO, Khoperskov AV, Khrapov SS. The problem of forecasting and control the hydrological regime in the mountainous area in the period flash floods on the basis of hydrodynamical numerical experiments. *Kibernetika i programirovanie*. 2016;**3**:35-53

Failure Analysis of a High-Pressure Natural Gas Heat Exchanger and its Modified Design

Lei-Yong Jiang, Yinghua Han, Michele Capurro and
Mike Benner

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.71202>

Abstract

The beauty of numerical simulations is its ability to reveal the physics or nature of practical engineering problems in detail, and then, to identify adequate solutions. In this chapter, an excellent example is demonstrated. The rupture of a heavy-duty, high-pressure natural gas heat exchanger is numerically investigated, and the importance of gravity effect is identified, which is often considered as a trivial factor. For the original design, the natural convection in the flow field of the heat exchanger is comparable with the forced convection at the designed operating conditions. These two convections are perpendicular and compete with each other, the flow field is highly unsteady, and high-temperature natural gas is trapped in the upper portion of the vessel, which causes the damage of the exchanger. By vertically mounting the exchanger assembly and locating the outlet pipe on top of the exchanger, the flow parameters become rather uniform at each vertical cross section and the wall temperature of the heat exchanger remains more or less the same as the heated natural gas. The proposed design has been successfully used up to now.

Keywords: forced convection, natural convection, heat transfer, heat exchanger, turbulent diffusion

1. Introduction

The beauty of numerical simulations is its ability to reveal the detailed phenomena or nature of complicated practical engineering problems, which are difficult and sometimes impossible from experimental studies. Based on the obtained numerical results, adequate solutions can easily be identified.

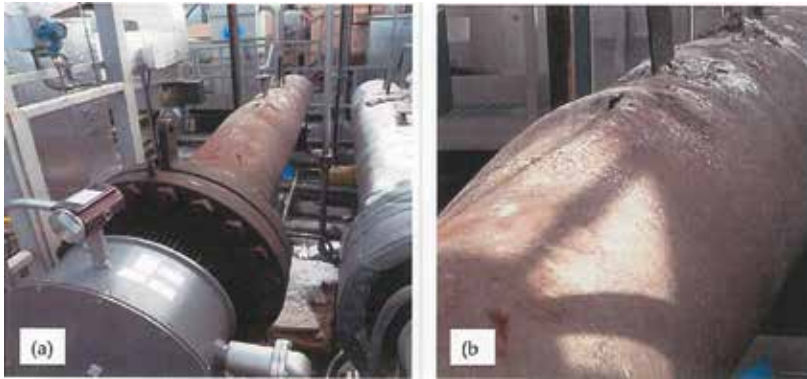


Figure 1. The natural gas heat exchanger: (a) the damaged heat exchanger; and (b) a close view of the cracks on the top surface of the vessel.

This chapter gives an excellent example for this type of approaches. To identify the reasons for the rupture of a heavy-duty, high-pressure natural gas heat exchanger, as shown in **Figure 1**, the flow field of the heat exchanger was numerically examined. Based on the findings, a new configuration was suggested and the corresponding flow field was studied. The installed modified heat exchanger has been trouble-free used to date. In the following sections, the computational domain, mesh, numerical methods, flow features of the two designs, and the cause of the heat exchanger damage are presented and discussed.

2. Numerical simulations

2.1. Computational domain and mesh of the original design

The natural gas heat exchanger has a nominal power of 390 KW, and the effective length of its heavy-duty vessel is 3785 mm with an outside diameter of 457.2 mm. As a sketch shown in **Figure 2**, it accommodates 138 (276 rods) heating elements. The diameter of these elements is 10.9 mm, and the length is 3277 mm for 64 long elements and 3252 mm for 74 short elements.

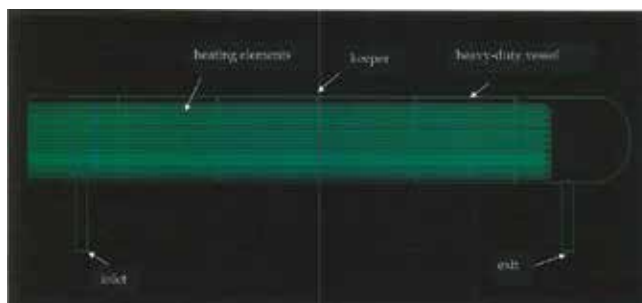


Figure 2. The computational domain of the original design.

Five keepers are inserted inside the vessel to maintain the proper radial positions of these elements. The nominal diameter of the inlet and exit pipes is 80 mm. The computational domain covers the whole flow field of the heater from the inlet to the exit, including heat elements and five keepers. It is important to mention that the whole natural gas heat exchanger assembly was mounted horizontally.

The mesh for one section of the heat exchanger vessel is shown in **Figure 3(a)**, and the mesh at a cross section cutting through heating elements is illustrated in **Figure 3(b)**. In these two plots, the meshed areas are where the natural gas flows, and the unmeshed hollow regions or circles are where the heating elements are located. **Figure 3(c)** is the mesh cutting through one keeper. As shown in **Figure 3(c)**, there are hundreds of small holes (11.5 mm in diameter) on the perforated plate of the keeper, 276 holes are considered blocked by the heating elements, and the rest meshed are flow passages. A narrow annular flow channel surrounding the keeper is used to keep the heating elements away from the vessel inner surfaces. Adjacent to the annular flow channel, parts of full circles are cut out by a flat bar of the keeper. The mesh size is ~ 4.0 million in the number of cells.

2.2. Boundary conditions

The inlet boundary conditions for the numerical simulations are listed in **Table 1**.

The heavy-duty vessel was wrapped with insulation material, and so it was reasonable to assume adiabatic boundary condition for its external wall boundaries. A heat flux of 4.48 KW/m^2 was specified at the heating-element surfaces starting from the middle cross section of the inlet pipe to the end of the heater elements. An increase in natural gas temperature by $\sim 200 \text{ K}$, from the room temperature, was expected. In the investigation, the natural gas was considered as pure methane.

2.3. Numerical methods

Steady, turbulent, thermal flows were considered in the present work, and a commercial software package, Fluent, was used for all simulations. The governing Favre-averaged conservation equations of mass, momentum, and enthalpy are not reproduced here, but can be readily found in [1, 2].

For closure of these partial differential equations, the realizable $k-\varepsilon$ turbulence model was applied to model turbulent momentum transfer. A benchmark study on turbulence models indicated that this model was superior to other four popular two-equation models and could provide similar results as those from the Reynolds stress model, a second-momentum closure [3]. The Reynolds analogy [4] was used to account for turbulent enthalpy transfer, and for this type of pipe flows, the turbulence Prandtl number of 0.7 was used [5, 6]. The gravity of 9.8 m/s^2 was assigned in the direction consistent with the heat exchanger mounting orientation.

For the thermal properties of methane, polynomials derived from the NIST JANAF tables [7] were used to calculate the specific heat as a function of temperature. Data from NIST [8] were used to obtain polynomials to determine the molecular viscosity and thermal conductivity of methane as functions of temperature.

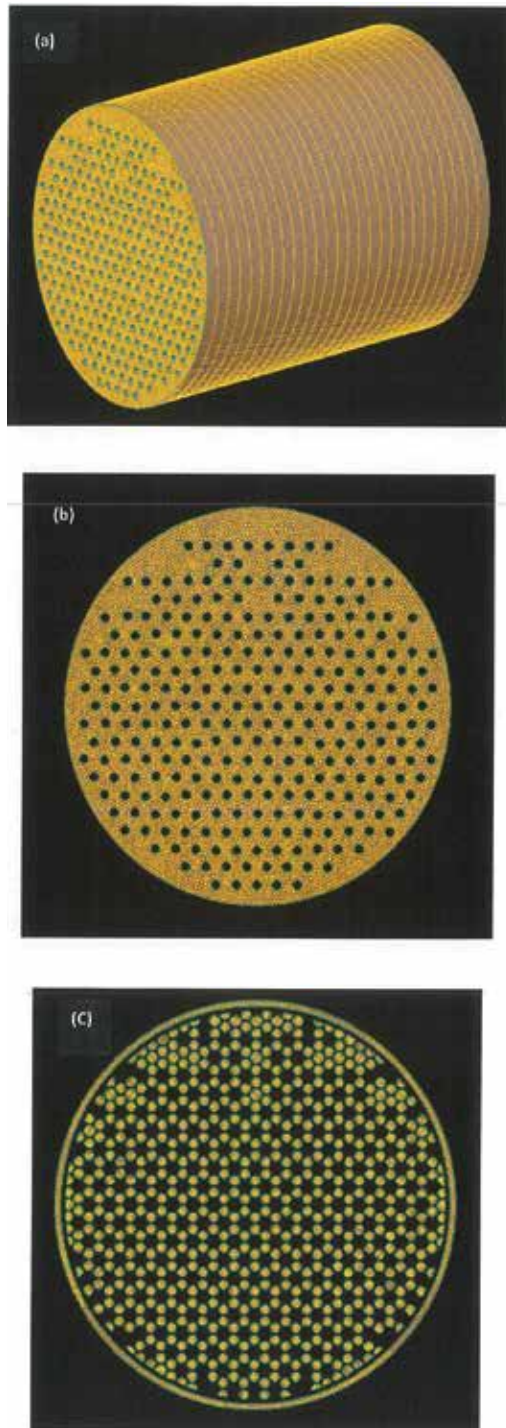


Figure 3. Meshes: (a) mesh for one section of the heat exchanger; (b) mesh at the section across heating elements; and (c) mesh across one heating-element keeper.

Mass flow rate	0.237 kg/s
Temperature	295.3 K
Absolute pressure	30.4 bar

Table 1. The inlet boundary conditions.

A segregated solver with a second-order accuracy scheme was chosen to resolve the flow fields. At convergence, the imbalance of mass flow rate between the inlet and exit was less than 0.34% for the original design and 0.007% for the modified configuration, while for the energy imbalance it was 0.38% for the former and 0.007% for the later. Due to the unsteady nature of the thermal flow field of the original design, the convergence could not reach the level for the modified case. Sixteen cores of a 64-bit LINUX cluster with 4 GB RAM for each core were used to perform all simulations.

3. Results and discussion

3.1. Results of the original design

Figure 4 shows the temperature contours along the longitudinal symmetric plane of the heat exchanger. Significantly, variation of temperature inside the heater is observed. High-temperature region exists in the upper portion of the vessel, while the low-temperature regions occur in the lower portion. The temperature profiles along the top and bottom walls are displayed in **Figure 5**. It is obvious that the top wall temperature reaching ~1700 K is considerably higher than the allowed service temperature of steel pipes, SA-106 GR.B [9]. Certainly, the heavy-duty vessel could not survive at such high-temperature and high-pressure operating conditions. Another important feature in **Figure 4** is the unsteady nature of the thermal flow field. Large eddies or fluid pockets randomly occur in the regions along and above the central axis of the vessel.

Why is the vessel wall temperature so high when the natural gas is set to be heated by only ~200 K? And why is the thermal flow field so unsteady in nature? These questions can be answered by analyzing the flow features or characteristics inside the heat exchanger vessel.

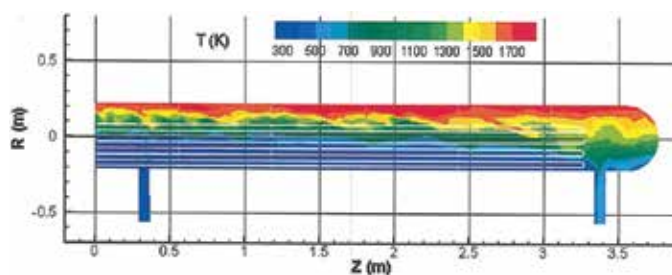


Figure 4. Temperature contours at the longitudinal symmetric plane.

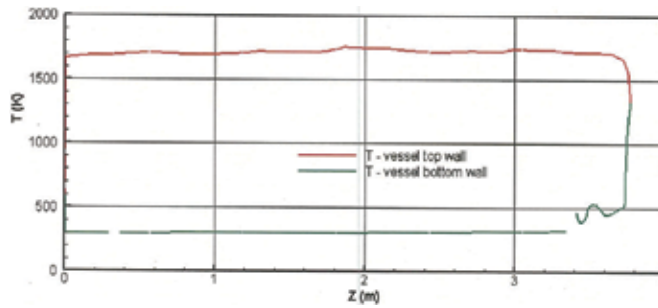


Figure 5. Temperature profiles along the vessel top and bottom walls.

The first or primary cause for the heat exchanger damage is that the forced convection in the flow field is weak, the natural convection is strong and comparable with the former, and these two actions are perpendicular and compete with each other. **Figure 6** is the velocity magnitude contour plot across the symmetric plane. As observed, the velocity inside the heater vessel is low, the mean velocity magnitude over the whole domain is only 0.174 m/s, and the mean Reynolds number is ~ 2800 . This means that the forced convection mainly in the longitudinal direction is weak. The absolute pressure distribution in the vessel is shown in **Figure 7**, and it varies a little around 3.04×10^6 Pa.

The methane density contours at the symmetric plane are shown in **Figure 8**. The density changes dramatically inside the heat exchanger vessel. High-density regions appear in the inlet pipe and lower portion of the vessel, and the maximum value reaches 20 kg/m^3 , while the low-density regions happen in the upper portion of the vessel with a minimum of 3.0 kg/m^3 . Due to the gravity, large differences in density induce strong natural convection inside the vessel.

As shown in **Figure 4**, the unsteady flow feature, randomly distributed flow pockets, is also observed in **Figure 8**. Velocity vectors at the portion of the symmetric plane are illustrated in **Figure 9**, where the vessel and keeper walls are indicated by blue lines, the length of vectors represents the magnitude of local velocities, and for comparison a reference vector of 0.2 m/s is provided. As shown in **Figure 9**, large counterclockwise recirculation regions are formed in the upper half of the vessel, which are induced by the relatively high horizontal velocities in the

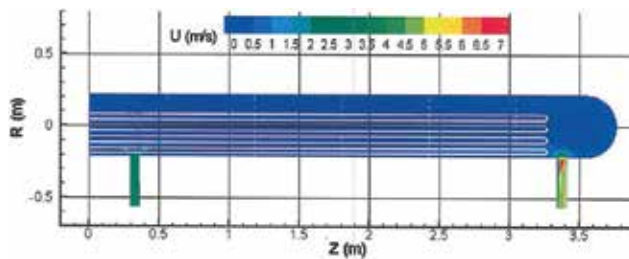


Figure 6. Velocity contours at the longitudinal symmetric plane.

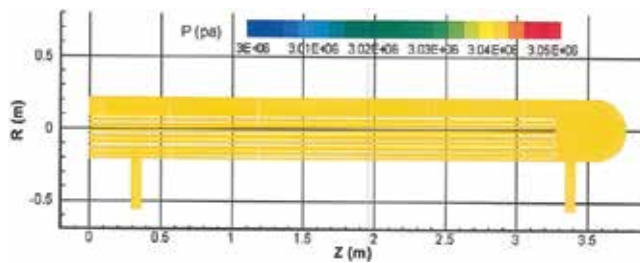


Figure 7. Absolute pressure contours at the longitudinal symmetric plane.

lower half of the vessel. The vertical velocity component corresponding to the natural convection is comparable with the horizontal component related to the forced convection. As a result of the competition between the two convections, the flow field inside the vessel becomes unstable in nature. The gas temperature at one point in the high-temperature region can vary by ± 80 K.

These observations are consistent with a first-order analytical assessment in Ref. [10]. As stated in the book, when the ratio of $R = Gr/(Re)^2 \approx 1$, the combined forced and natural convections must be considered in heat transfer analysis. Here, Gr is the Grashof number and Re stands for Reynolds number. Based on the averaged values of flow parameters, this ratio equals to 0.6 for this case.

Figures 10–13 provide detailed distributions of temperature and velocity magnitude at the inlet, outlet, and four middle cross sections, and at the five keeper cross sections, respectively. These plots further confirm the above observed flow features. The high-temperature region occupies about one-third of the cross-sectional areas, and the low-temperature region gradually decreases from half of the area at the inlet section, to about one-third of the area at the last keeper section, and eventually a fraction at the outlet section (**Figures 10 and 12**). As observed in **Figures 11 and 13**, the velocity magnitude in the lower halves at these cross sections is higher than that in the upper halves.

The second reason for the heat exchanger damage is that the flow exit is located at the lowest position of the whole flow domain (**Figure 2**). Consequently, the high-temperature or low-density fluid is trapped in the upper portion of the vessel, does not flow out of the vessel, and keeps recirculating, as shown in **Figures 9 and 14**.

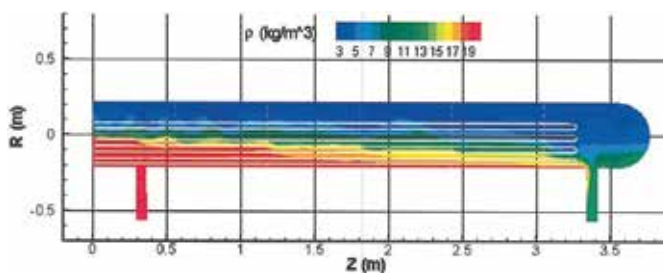


Figure 8. Density contours at the longitudinal symmetric plane.

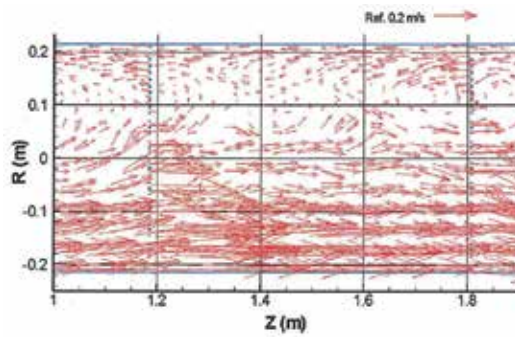


Figure 9. Velocity vectors at part of the longitudinal symmetric plane.

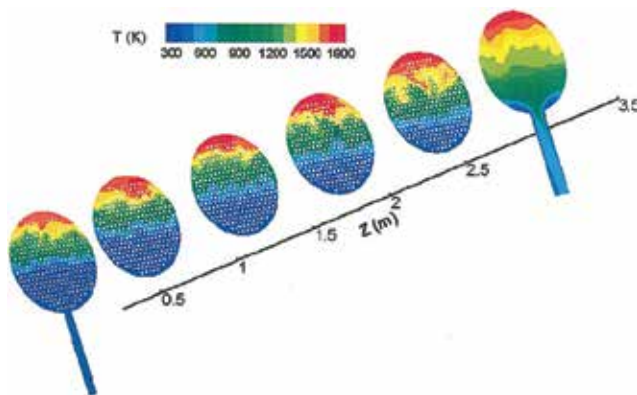


Figure 10. Temperature contours at the inlet, outlet, and four middle cross sections.

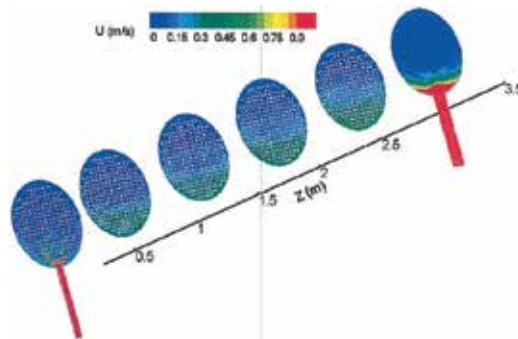


Figure 11. Velocity magnitude contours at the inlet, outlet, and four middle cross sections.

Notice that the fluid in the upper high-temperature regions is continuously heated by the heating elements that are more or less uniformly distributed over the vessel cross sections. The only way, for the fluid in these swirling regions to release some of the heat, is through diffusion (molecular and weak turbulent), which is significantly less effective than convection.

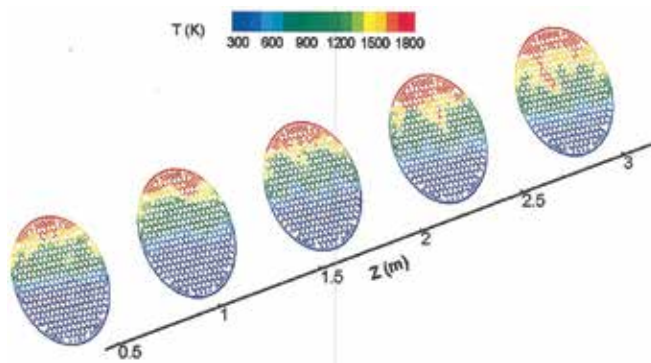


Figure 12. Temperature contours at the five keeper cross sections.

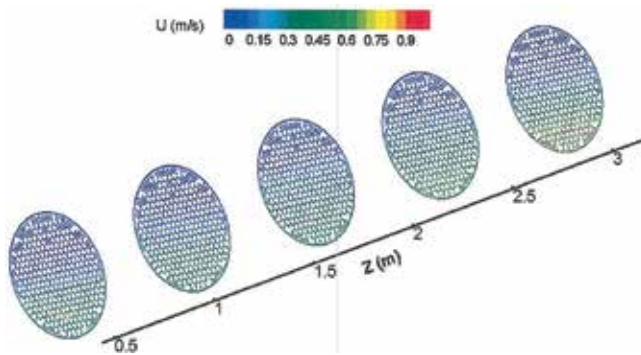


Figure 13. Velocity magnitude contours at the five keeper cross sections.

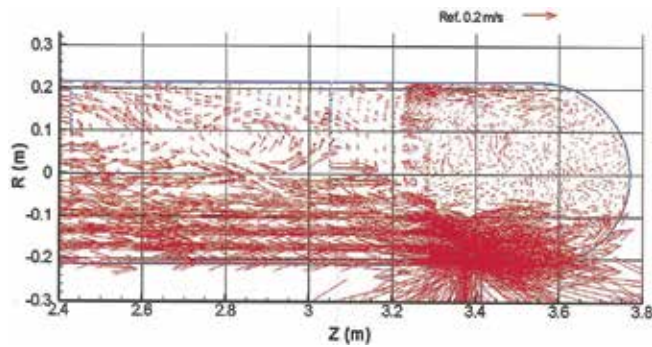


Figure 14. Velocity vectors at the downstream part of the longitudinal symmetric plane.

When the flow reaches quasi-steady, the gas temperature can be as high as ~1700 K (Figure 4). This is why although the mean methane temperature at the exchanger exit is increased by only ~200 K, the temperature at the top wall of the vessel can reach ~1700 K.

The above results and discussion suggest that to avoid the competition between the forced and natural convections in a perpendicular manner, the heat exchanger assembly should be mounted vertically, and to avoid fluid trapping, the flow exit pipe should be located on top of the vessel. With these arrangements, it is expected that the gravity effect or natural convection effect would be more or less uniform at each horizontal cross section, and no local high-temperature region would occur inside the heat exchanger.

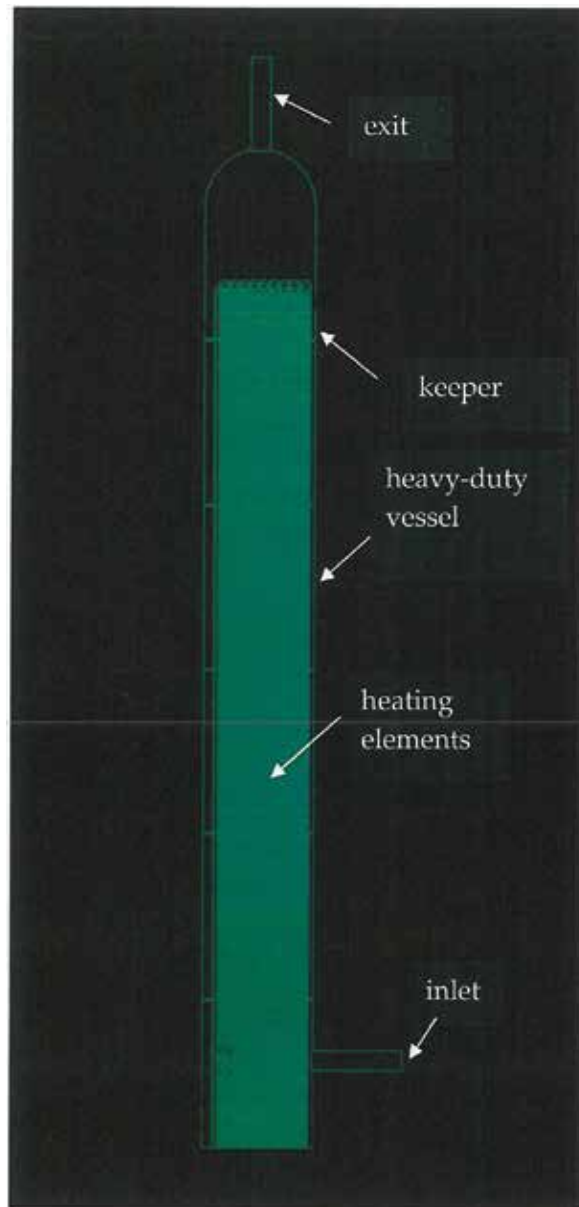


Figure 15. Modified heat exchanger configuration.

3.2. Results of the modified natural gas heat exchanger

The modified heat exchanger configuration is shown in **Figure 15**, where the whole assembly is mounted vertically and the outlet pipe is moved to the vessel top, and other parts remain the same as the original design. Similar mesh was generated for the new design, and the boundary conditions and numerical methods were unchanged.

The temperature contours along the longitudinal symmetric plane are shown in **Figure 16**. The flow temperature gradually increases from 295 K at the inlet to 495 K at the exit with an increase of 200 K. The maximum temperature is 564 K at the top surfaces of the heating elements

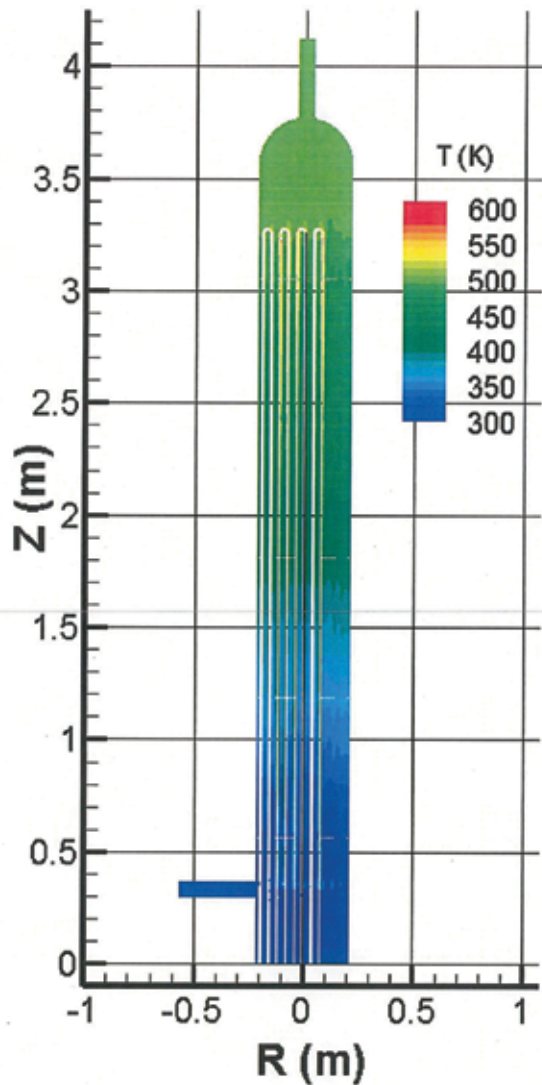


Figure 16. Temperature contours at the longitudinal symmetric plane.

(also see **Figure 21** later), and the temperature difference between the element walls and surrounding fluid is the driving force for heat transfer from the heat elements to the fluid. The temperature profiles at the right- and left-side walls are displayed in **Figure 17**. The wall temperature gradually increases along the vertical direction from 295 K to 495 K, and the maximum wall temperature is equal to the exit gas temperature.

Similar to the original design, the velocity magnitude shown in **Figure 18** is low inside the vessel, and its averaged value is 0.18 m/s with a maximum of 5.3 m/s at the center of the exit. The absolute pressure distribution inside the vessel also varies a little around 3.04×10^6 Pa,

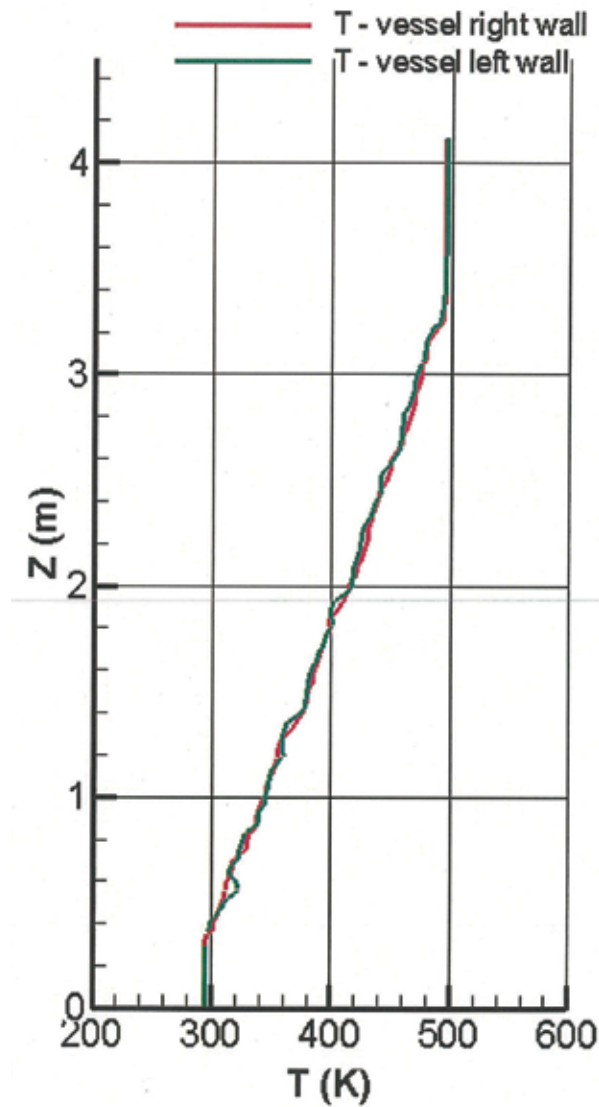


Figure 17. Temperature profiles along the vessel right and left walls.

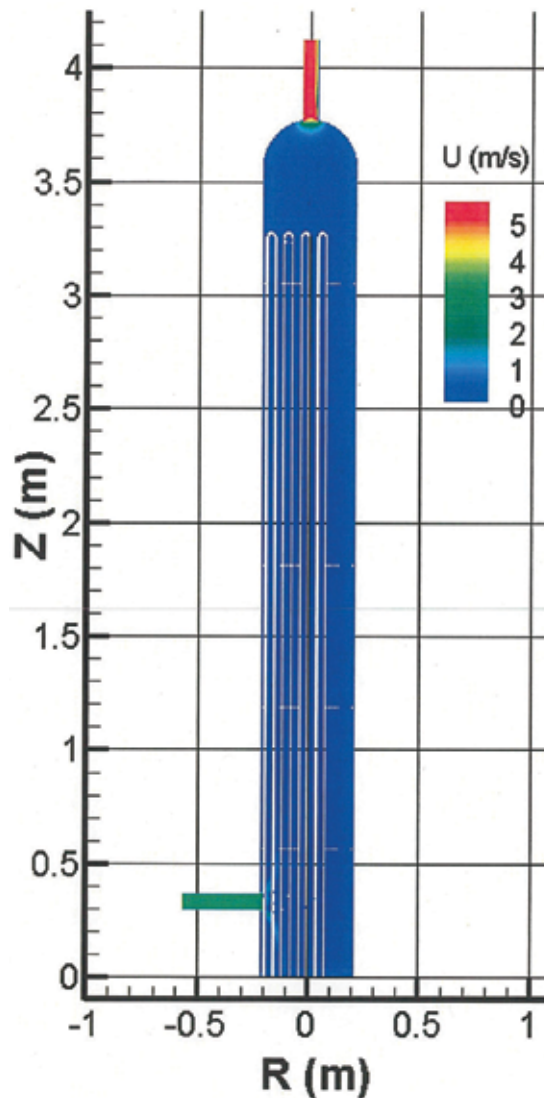


Figure 18. Velocity magnitude contours at the longitudinal symmetric plane.

as indicated in **Figure 19**. **Figure 20** presents the density contours at the symmetric plane. It gradually decreases from 20 kg/m^3 at the inlet to 11.8 kg/m^3 at the exit and is more or less uniform at vertical cross sections.

Detailed distributions of temperature and velocity magnitude at the inlet, five middle and exit cross sections are provided in **Figures 21** and **22**. The same parameter plots across the five keepers are given in **Figures 23** and **24**. These figures clearly indicate that the flow parameters are rather uniform at each cross section, particularly at the five keeper cross sections. The temperature gradually increases from the upstream to downstream sections, as illustrated in

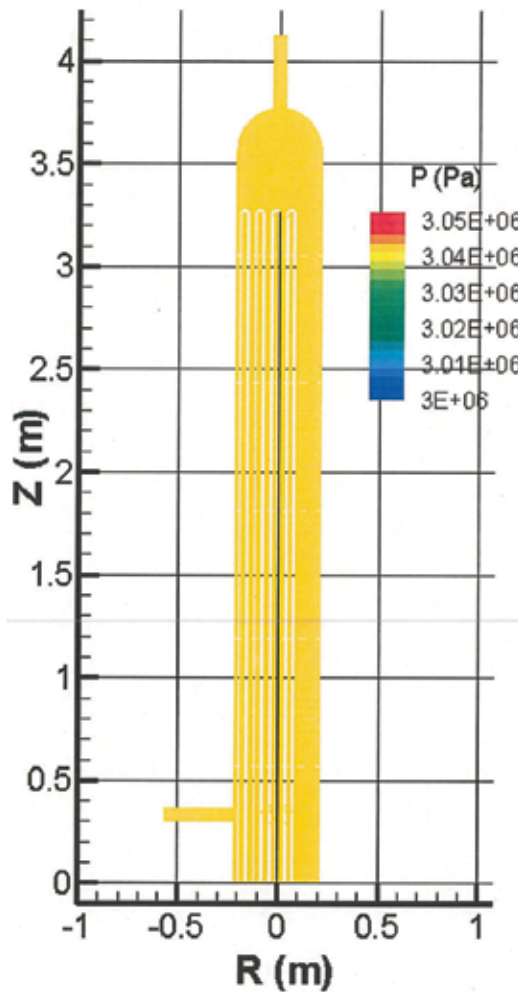


Figure 19. Absolute pressure contours at the longitudinal symmetric plane.

Figures 21 and 23, and the maximum temperature is about 560 K. As shown in Figure 22, the flow velocity gradually increases from the inlet to the exit, except for small local regions at the inlet and first middle sections.

In summary, the flow features of the modified design are remarkably different from those for the original design, except that the velocity magnitude is low and the absolute pressure is about 3.4×10^6 Pa for both cases. For the modified design, both the natural and forced convections are aligned in the vertical direction; therefore, the flow parameters are more or less uniform at each vertical cross section and the flow field is stable without randomly located large recirculation regions. Moreover, the gas flows out of the exit at the required temperature and the vessel wall temperature remains the same as the surrounding gas. This proposed design has been successfully used to date.

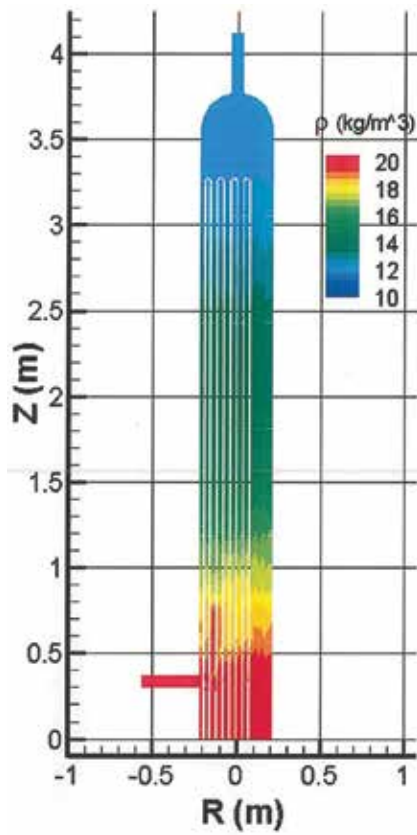


Figure 20. Density contours at the longitudinal symmetric plane.

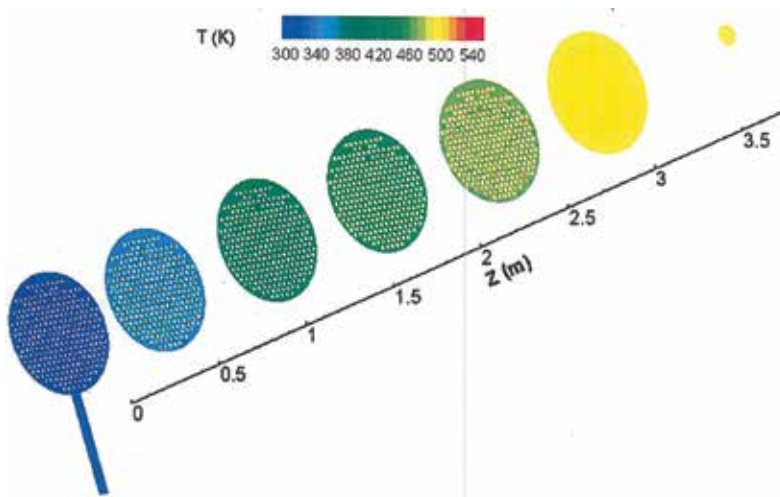


Figure 21. Temperature contours at the inlet, outlet, and five middle cross sections.

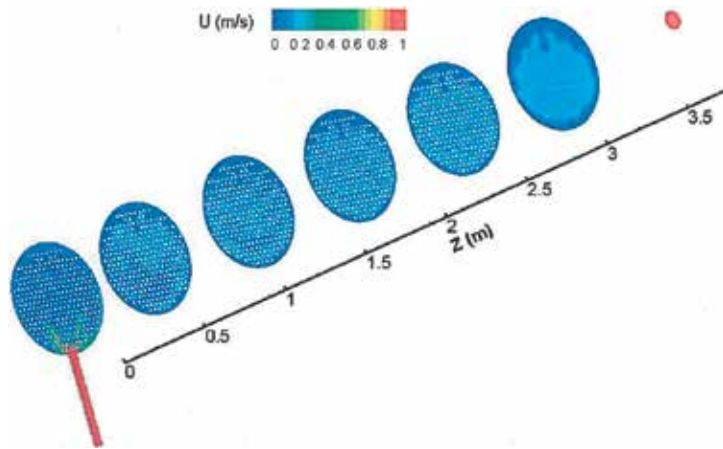


Figure 22. Velocity magnitude contours at the inlet, outlet, and five middle cross sections.

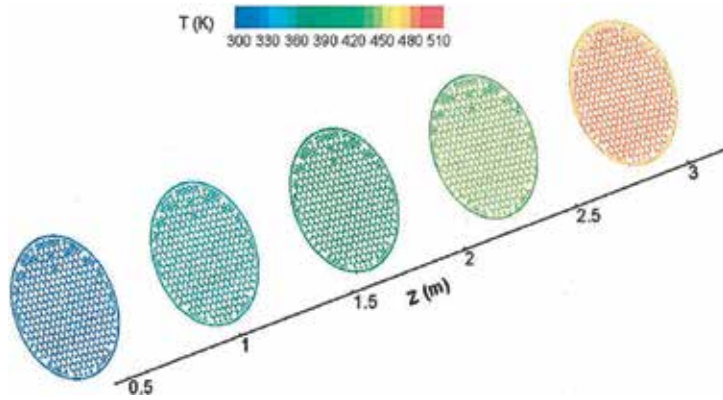


Figure 23. Temperature contours at the five keeper cross sections.

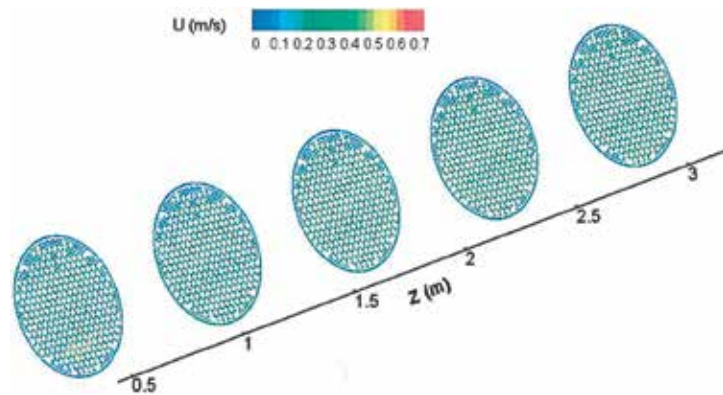


Figure 24. Velocity magnitude contours at the five keeper cross sections.

4. Conclusions

To investigate the damage of a natural gas heat exchanger, the numerical simulations of the flow fields of the original and modified designs are performed. It is found that there are two reasons for the damage. First, at the required operating conditions, the forced convection is weak, and the natural convection is strong and comparable with the forced convection. These two actions are perpendicular and compete to each other. As a result, strong unsteadiness in the flow field is induced. Second, the whole assembly is mounted horizontally and the flow exit pipe is located at the lowest position. Consequently, the high-temperature or low-density fluid is trapped in the upper portion of the vessel. The trapped fluid is continuously heated by the heating elements located in the upper region of the vessel and eventually exceeds the allowed service temperature of the steel pipe.

The numerical results and analysis suggest that the heat exchanger assembly should be mounted vertically and the exhaust pipe should be located at the top of the exchanger. With these modifications, the flow parameters become more or less uniform at each vertical cross section, the flow field becomes stable, the methane temperature at the exit reaches the designed value, and the vessel wall temperature remains the same as the surrounding gas. This new design has been trouble-free used up to now.

Author details

Lei-Yong Jiang*, Yinghua Han, Michele Capurro and Mike Benner

*Address all correspondence to: lei-yong.jiang@nrc-cnrc.gc.ca

Aerospace, National Research Council of Canada, Ottawa, Ontario, Canada

References

- [1] Poinso T, Veynante D. *Theoretical and Numerical Combustion*. Philadelphia, PA: R. T. Edwards Inc.; 2005
- [2] ANSYS Inc., *Fluent 18.0 Documentation*. Lebanon, NH, USA; 2016
- [3] Jiang LY. A critical evaluation of turbulence modelling in a model combustor. *ASME Journal of Thermal Science and Engineering Applications*. 2013;5(3):031002
- [4] Jiang LY, Campbell I. Reynolds analogy in combustor Modelling. *International Journal of Heat and Mass Transfer*. 2008;51(5-6):1251-1263
- [5] Hinze JO. *Turbulence*. New York: The McGraw-Hill Book Company Inc.; 1987. p. 372-753
- [6] White FM. *Heat and Mass Transfer*. New York: Addison-Wesley Publishing Company; 1988. p. 320-641

- [7] Chase MW. National Institute of Standards and Technology (U.S.), NIST-JANAF Thermochemical Tables, 4th edition, Washington, DC, 1998
- [8] Lemmon EW, McLinden MO, Huber ML. REFPROP, Reference Fluid Thermodynamic and Transport Properties, NIST Standard Reference Database 23, Version 7.0, 2002
- [9] OneSteel Piping Systems. 2012. Available from: http://www.onesteelbuildingservices.com/pdffiles/OneSteel_%20Pipe_Catalogue_web.pdf
- [10] Incropera FP, DeWitt DP. Fundamentals of Heat and Mass Transfer. USA: John Wiley & Sons Inc.; 2002

Numerical Simulation of Nanoparticles with Variable Viscosity over a Stretching Sheet

Noreen Sher Akbar, Dharmendra Tripathi and
Zafar Hayat Khan

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.71224>

Abstract

The effects of different types of base fluids on carbon nanotube (CNT) nanofluids flow over a circular stretching sheet are numerically analyzed. The nonlinear variation of radial velocity in radial direction is assumed at surface of stretching sheet. The temperature dependent fluid viscosity is taken into consideration. Two different types of flows (assisting flow and opposing flow) are discussed under the buoyant force effects. Single walled CNT and multi walled CNT are considered as nanoparticles for better thermal conductivity of the nanofluids. A set of similarity transformations to convert the partial differential equations into ordinary differential equations is hired. The non-linear ODEs are numerically solved by employing fourth order Runge-Kutta method. Discussions of numerical simulations for flow characteristics have been made appropriately. A comparative study for various type of base fluids like kerosene, engine oil and ethylene glycol is also presented. From the predicted simulation, it is observed that the variation in Nusselt number is maximum for engine oil and minimum for kerosene oil however, the variation in skin friction coefficient is largest for kerosene oil and least for engine oil. Furthermore, numerical results are also validated with achieving a good correlation with existing results.

Keywords: CNT nanofluids, nonlinear stretching sheet, Runge-Kutta method, buoyancy force, heat transfer, similarity transformation

1. Introduction

The first key point of present investigation is carbon nanotubes (CNTs) nanofluids that means the suspension of CNTs in conventional fluids such as air, helium, water, minerals oils, Freon, and ethylene glycol. The applications of nanofluids spread over a wide range of disciplines, including heat transfer, material science, physics, and chemical engineering. Many techniques have been used to enhance the thermal conductivity of the base fluids in which the suspension

of micro/nano-sized particles in fluids have also been tried since many decades. However, the word 'nanofluids' was primarily introduced by Choi [1]. He has studied the thermal conductive effects on convectational fluids subject to suspension of metallic nanoparticles. He has observed that with suspension of nanoparticles, the thermal conductivity of convectational fluids enhances. In these directions, numerous investigations have therefore been carried out in the past few decades, seeking to wide applications and developments, some interesting reviews [2–4] have been reported.

To study the heat transfer rate of water based nanofluids, the various types of nanoparticles like titanium dioxide, alumina, silica diamond, zinc-oxide and copper [5–7] have been considered. It has been depicted that the thermal conductivity enhances with suspension of the nanoparticles. Moreover, CNTs have also received great interest due to significant enhancement of thermal conductivity, unique structure and physical (mechanical and electrical) properties [8, 9]. Ding et al. [8] have prepared a nanofluids with suspension of CNTs in distilled water and measured the thermal conductivity and viscosity of CNTs nanofluids. They have reported the enhancement of thermal conductivity depends on CNTs concentration, and pH level. They also concluded that nanofluids with 0.5 wt.% CNTs, the maximum enhancement is over 350% at $Re = 800$. Ko et al. [9] experimentally reported the flow characteristics of CNTs nanofluids in a tube and summarized that the friction factor for CNT nanofluid is low in compare to water. Another experimental study for thermophysical properties of CNTs nanofluids with base fluid as mixture of water and ethylene glycol has been presented by Kumaresan and Velraj [10]. And this study noted that the maximum thermal conductivity enhances up to 19.75% for the nanofluid containing 0.45 vol.% MWCNT at 40°C. In addition, few more experimental investigations [11, 12] with applications in a tubular heat exchanger of various lengths for energy efficient cooling/heating system and turbulent flow heat exchanger are performed. With CNTs nanofluids, some more numerical simulations [13, 14] for flow characteristics are presented in literature and discussed the physiological flows application.

The second key point is boundary layer flow over a circular stretching sheet. The boundary layer flow past a stretching sheet was initially investigated by Crane [15]. He has discussed its applications such as annealing and tinning of copper wires, polymer extrusion in melt spinning process, manufacturing of metallic sheets, paper production, and glass, fiber and plastic production etc. The heat transfer rate play important role in all manufacturing, productions and fabrication processes. This model has been explored by many investigators for various physical aspects. However, some interesting extensions of Crane's model for boundary layer flow of nanofluids have recently investigated, see examples [16–30]. The first boundary layer flow model for nanofluids flow past stretching sheet was presented by Khan and Pop [16] which was extended for a convective boundary condition [17], nonlinear stretching sheet [18], unsteady stretching surface [19], micropolar nanofluid flow [20], magneto-convective non-Newtonian nanofluid slip flow over permeable stretching sheet [21], Non-aligned MHD stagnation point flow of variable viscosity nanofluids [22], Stagnation electrical MHD mixed convection [23], exponential temperature-dependent viscosity and buoyancy effects [24], thermo-diffusion and thermal radiation effects on Williamson nanofluid flow [25], magnetic dipole and radiation effects on viscous ferrofluid flow [26], transient ferromagnetic liquid flow [27], magnetohydrodynamic Oldroyd-B nanofluid [28],

spherical and non-spherical nanoparticles effects [29], three dimensional free convective magnetohydrodynamics [30]. Moreover, some works [31–34] extended Khan and Pop’s model for CNTs nanofluids where combined effects of slip and convective boundary conditions have been discussed [31], convective heat transfer in MHD slip flow has been studied [32], nonlinear stretching sheet with variable thickness has been considered [33] and variable thermal conductivity and thermal radiation effects have been discussed.

After reviewed many investigations [16–34] on boundary layer flow of nanofluids past through linear/nonlinear stretching sheet, none of the them has been reported for boundary layer flow of CNTs nanofluids over a circular stretching sheet. Considering this gap of research, we formulate a model and also numerically simulate it to study the flow characteristic of CNTs nanofluids flow over a circular stretching sheet. Kerosene, Engine oil and Ethylene glycol are considered for base fluids. And MCNT and SWCNT are also taken as nanoparticles. The negative and positive buoyant force effects as opposing and assisting flows are also examined. The effects of fluid and flow parameters on velocity profile, temperature profile, skin-friction coefficient, Nusselt number and stream lines are computed numerically.

2. Mathematical model

We consider the laminar incompressible flow of CNTs nanofluids over a circular sheet aligned with the $r\theta$ -plane in the cylindrical coordinate system (r, θ, z) . The schematic representation and coordinate system are depicted in **Figure 1**.

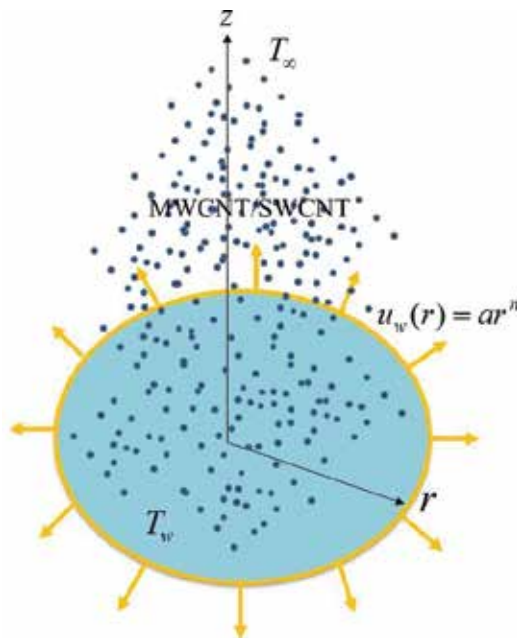


Figure 1. Schematic representation of CNTs nanofluids flow over circular stretching sheet.

The flow regime is considered as in the half space $z \geq 0$ and the sheet is stretched in radial direction with polynomial variation. The temperature at sheet is fixed as T_w and ambient temperature is considered as T_∞ . Considering the boundary layer flow assumptions, the governing equations for conservation of mass, momentum, and energy can be expressed as:

$$\frac{\partial u}{\partial r} + \frac{u}{r} + \frac{\partial w}{\partial z} = 0, \quad (1)$$

$$\left(u \frac{\partial u}{\partial r} + w \frac{\partial u}{\partial z} \right) = \frac{1}{\rho_{nf}} \frac{\partial}{\partial z} \left(\mu_{nf}(T) \frac{\partial u}{\partial z} \right) \pm g\beta(T - T_\infty), \quad (2)$$

$$\left(u \frac{\partial T}{\partial r} + w \frac{\partial T}{\partial z} \right) = \alpha_{nf} \frac{\partial^2 T}{\partial z^2}, \quad (3)$$

where u and w are the velocity components along r - and z -directions respectively, T is the temperature. Here the “+” sign in Eq. (2) corresponds to the assisting flow while the “-” sign corresponds to the opposing flow, respectively. Here, ρ_{nf} is the nanofluid density, k_f is the thermal conductivity of the fluid, T_w is the wall temperature, $\mu_{nf}(T)$ is the temperature dependent viscosity of nanofluid and α_{nf} is the thermal diffusivity of nanofluid which are defined [31–34] as:

$$\mu_{nf} = \frac{\mu_0 e^{-\alpha\theta}}{(1 - \phi)^{2.5}}, \quad (4a)$$

$$\alpha_{nf} = \frac{k_{nf}}{(\rho c_p)_{nf}}, \quad \rho_{nf} = (1 - \phi)\rho_f + \phi\rho_s, \quad (4b)$$

$$(\rho c_p)_{nf} = (1 - \phi)(\rho c_p)_f + \phi(\rho c_p)_s, \quad (4c)$$

$$(\rho\gamma)_{nf} = (1 - \phi)(\rho\gamma)_f + \phi(\rho\gamma)_s, \quad (4d)$$

$$k_{nf} = k_f \left(\frac{k_s + 2k_f - 2\phi(k_f - k_s)}{k_s + 2k_f + 2\phi(k_f - k_s)} \right). \quad (4e)$$

Here, ρ_f is density of the base fluid, ρ_s is density of the nanoparticles, k_f is thermal conductivity of the base fluid, k_s is thermal conductivity of the nanoparticles, γ_{nf} is the thermal expansion coefficient, γ_f is the thermal expansion coefficient of base fluid, ϕ is the nanoparticle volume fraction, and γ_s is the thermal expansion coefficient of the nanoparticles.

The following boundary conditions are to be employed:

$$u = u_w(r) = ar^n, \quad T = T_w, \quad \text{at } z = 0, \quad (5a)$$

$$u \rightarrow 0, \quad T \rightarrow T_\infty, \quad \text{as } z \rightarrow \infty, \quad (5b)$$

in which u_w is the wall velocity, $a > 0$ is the stretching constant and $n > 0$ is the power-law index.

Introducing the following similarity transformations:

$$\eta = \sqrt{\frac{a}{v_f}} r^{\frac{(n-1)}{2}} z, u = ar^n f'(\eta), \theta = \frac{T - T_\infty}{T_w - T_\infty}, w = -ar^{\frac{(n-1)}{2}} \sqrt{\frac{v_f}{a}} \left(\frac{(n+3)}{2} f(\eta) + \frac{(n-1)}{2} \eta f'(\eta) \right). \quad (6)$$

Reynolds model of viscosity expression can be taken as:

$$\mu_f(\theta) = e^{-(\alpha\theta)} = 1 - (\alpha\theta) + O(\alpha^2), \quad (7)$$

where α is the viscosity parameter.

Using Eqs. (6) and (7) in Eqs. (1)–(5), the governing equations and boundary conditions are transformed as:

$$\frac{1 - (\alpha\theta)}{(1 - \phi)^{2.5}} f''' + \frac{(-\alpha\theta')}{(1 - \phi)^{2.5}} f''' + \left(1 - \phi + \phi \frac{\rho_s}{\rho_f} \right) \left\{ \frac{(n+3)}{2} f f'' - n (f')^2 \right\} \pm G_r \theta = 0, \quad (8)$$

$$\left(\frac{k_{nf}}{k_f} \right) \theta'' + \text{Pr} \left(\frac{n+3}{2} \right) \left(1 - \phi + \phi \frac{(\rho c_p)_s}{(\rho c_p)_f} \right) [(f\theta')] = 0, \quad (9)$$

$$f(0) = 0, f'(0) = 1, f'(\infty) = 0, \theta(0) = 1, \theta(\infty) = 0, \quad (10)$$

where $\text{Pr} = \frac{(\mu_0 c_p)_f}{k_f}$ is the Prandtl number.

The skin friction coefficients (C_{fx}) is defined as:

$$C_{fx} = \frac{\mu_{nf}(T)}{\rho_f u_w^2} \left(\frac{\partial u}{\partial z} \right)_{z=0}. \quad (11)$$

And the local Nusselt number (Nu_x) is defined as:

$$Nu_x = \frac{-x k_{nf}}{k_f (T_w - T_\infty)} \left(\frac{\partial T}{\partial z} \right)_{z=0}. \quad (12)$$

The dimensionless form of skin friction coefficients can be obtained as:

$$(\text{Re}_x)^{1/2} C_{fx} = \frac{\mu_f(\theta(0)) f''(0)}{(1 - \phi)^{2.5}}. \quad (13)$$

And the dimensionless form of skin friction coefficients can be obtained as:

$$(\text{Re}_x)^{1/2} Nu_x = -\frac{k_{nf}}{k_f} \theta'(0). \quad (14)$$

3. Numerical scheme

Numerical solutions of ordinary differential Eqs. (8) and (9) subject to boundary conditions (10) are obtained using a shooting method. First we have converted the boundary value problem (BVP) into initial value problem (IVP) and assumed a suitable finite value for the far field boundary condition, i.e. $\eta \rightarrow \infty$, say η_∞ . To solve the IVP, the values for $f''(0)$ and $\theta'(0)$ are needed but no such values are given prior to the computation. The initial guess values of $f''(0)$ and $\theta'(0)$ are chosen and fourth order Runge-Kutta method is applied to obtain a solution. We compared the calculated values of $f(\eta)$ and $\theta(\eta)$ at the far field boundary condition $\eta_\infty (=20)$ with the given boundary conditions (10) and the values of $f''(0)$ and $\theta'(0)$, are adjusted using Secant method for better approximation. The step-size is taken as $\Delta\eta=0.01$ and accuracy to the fifth decimal place as the criterion of convergence. It is important to note that the dual solutions are obtained by setting two different initial guesses for the values of $f''(0)$.

4. Numerical simulations and discussion

This section presents the numerical simulations for dimensionless velocity, temperature, skin friction, Nusselt numbers and streamlines under the influence of the flow parameters which are illustrated through the (Figures 2–11). Table 1 is also given for thermophysical properties of base fluids and CNTs (MWCNT and SWCNT).

Figures 2 and 3 illustrate the velocity profiles for assisting and opposing flows where multi walled CNT is considered as nanoparticle and Kerosene is taken as base fluid. The Grashof number is fixed at 0.5. From all the figures, it is depicted that $f'(\eta)$ attains maximum values at $\eta=0$ i.e. near the stretching sheet. It is clear from the transformation $\eta = \sqrt{\frac{g_c}{\nu_f}} r^{\frac{(n-1)}{2}} z$ that $\eta=0$

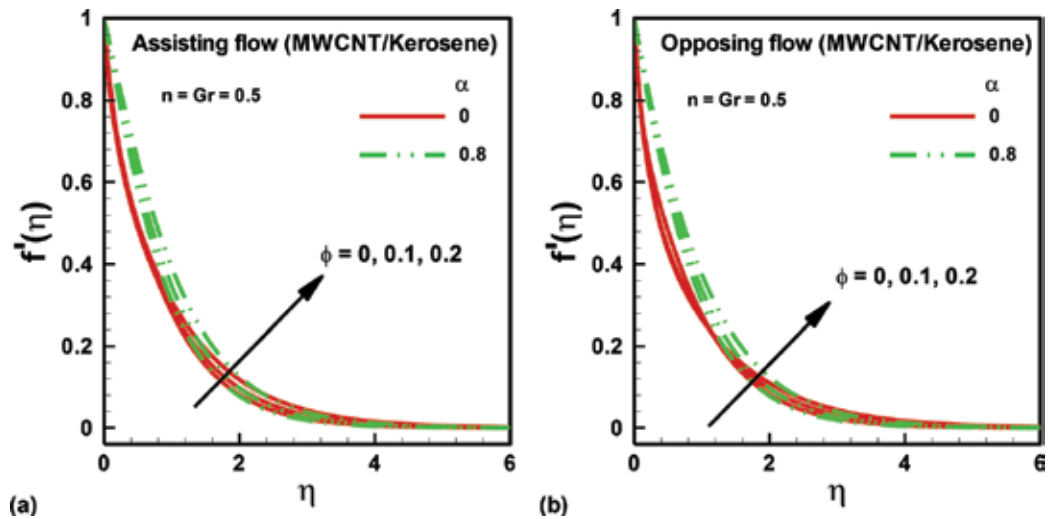


Figure 2. Velocity profile for different values of nanoparticle volume fraction and viscosity parameter: (a) assisting flow and (b) opposing flow.

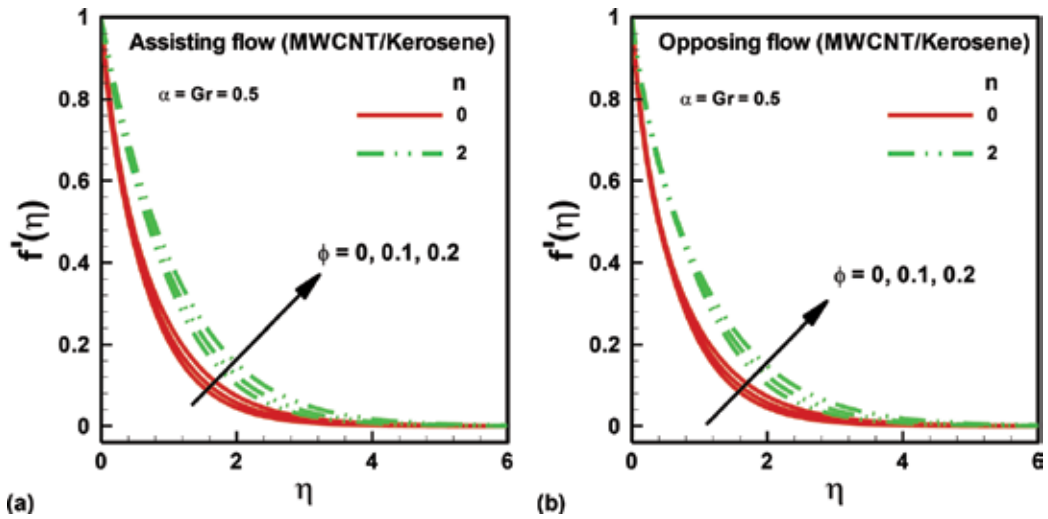


Figure 3. Velocity profile for different values of nanoparticle volume fraction and power law index: (a) assisting flow and (b) opposing flow.

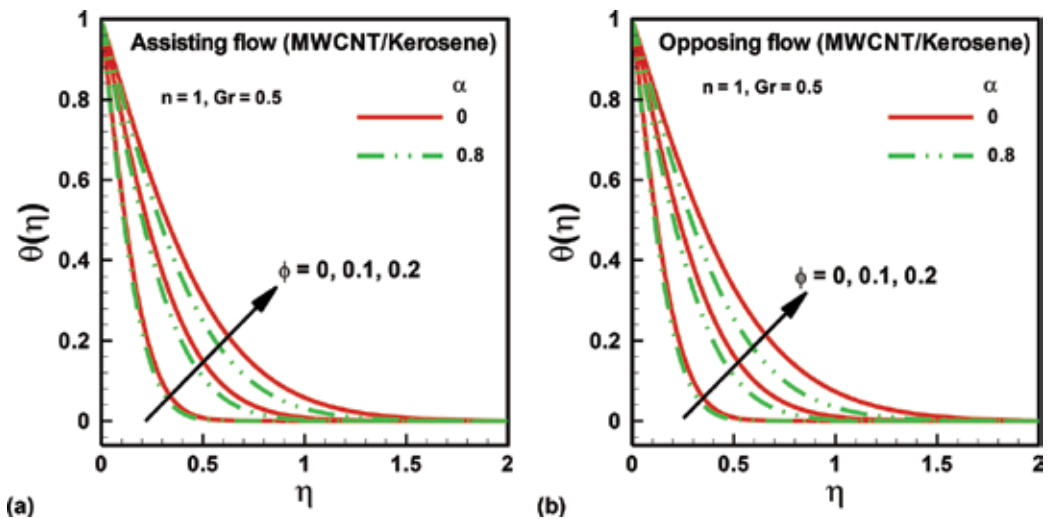


Figure 4. Temperature profile for different values of nanoparticle volume fraction and viscosity parameter: (a) assisting flow and (b) opposing flow.

at $z=0$. The curves show asymptotic nature long the η -axis. It is physically interpreted that the radial velocity diminishes when fluid is flowing away from the sheet. In **Figure 2 (a and b)**, the effects of nanoparticle volume fraction (ϕ) and viscosity parameter (α) on the velocity profiles for assisting flow and opposing flows are depicted and the value of power index is fixed at 0.5. It is inferred that with increasing the volume fraction parameter, the velocity slightly enhances for $\eta < 4$ for both types of flows. It is also revealed that with increasing the magnitude of viscosity parameter, the radial velocity goes up $\eta < 2$ for both types of flows. This is very clear from equation, $\mu_f(\theta) = e^{-(\alpha\theta)} = 1 - (\alpha\theta) + O(\alpha^2)$, that with increasing the magnitude of α , the

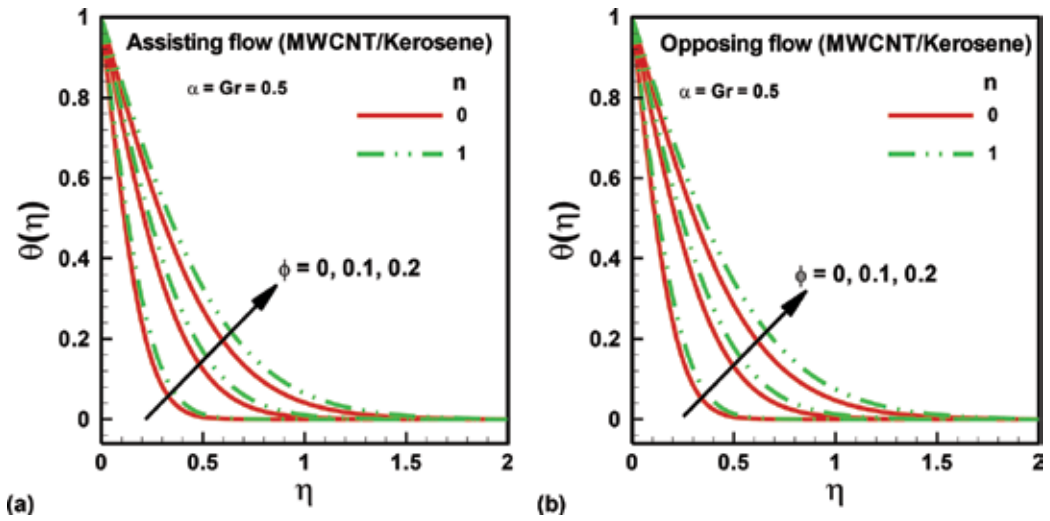


Figure 5. Temperature profile for different values of nanoparticle volume fraction and power law index: (a) assisting flow and (b) opposing flow.

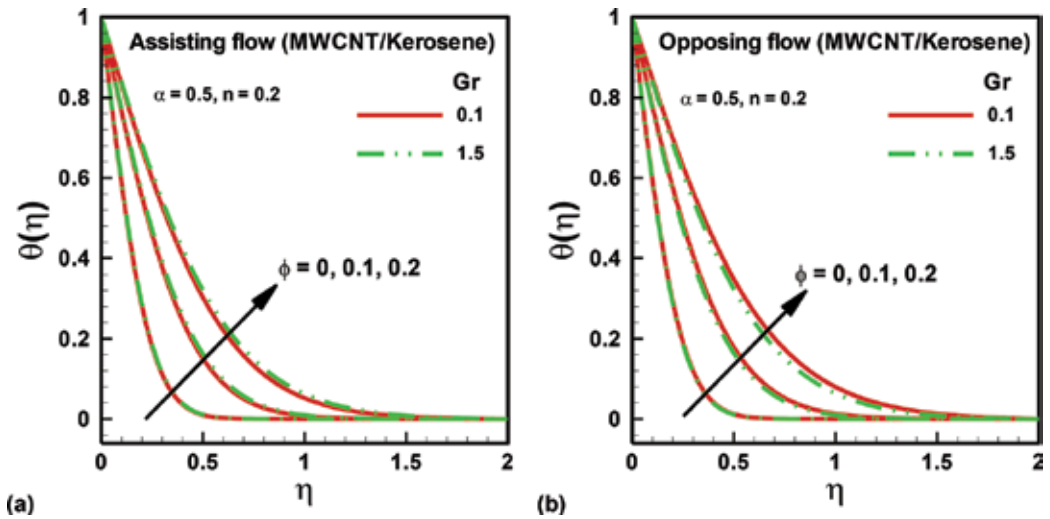


Figure 6. Temperature profile for different values of nanoparticle volume fraction and Grashof number: (a) assisting flow and (b) opposing flow.

viscosity will reduce. It is obvious that with increasing the viscosity of fluids, the shearing resistance will enhance and then velocity will go down. **Figure 3(a, b)** show that the effects of power law index (n) on velocity profile for both types of flows at fixed value of viscosity parameter 0.5. It is observed that the radial velocity rises with increasing the magnitude of power law index parameter. It is further noted that the velocity profile enlarges with increasing the nanoparticle volume fraction for all values of power law index and also viscosity parameter for both types of flows.

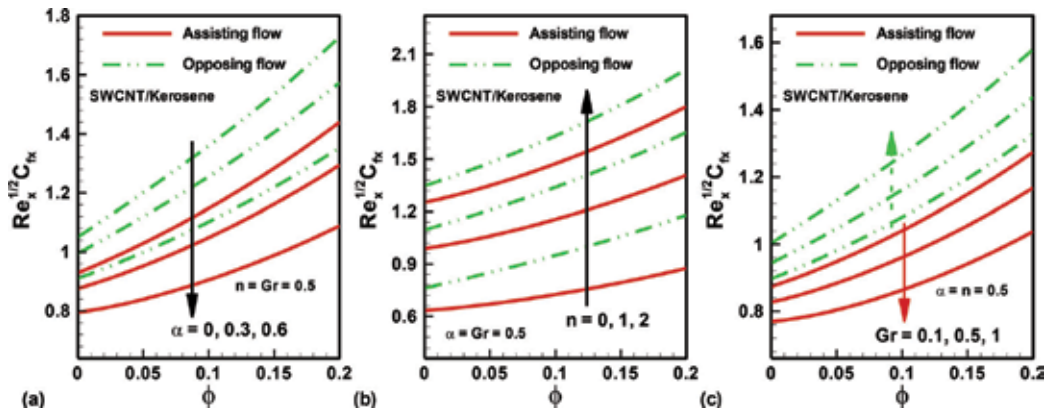


Figure 7. Variation of skin-friction coefficient against nanoparticle volume fraction for different values of (a) viscosity parameter, (b) power law index, and (c) Grashof number.

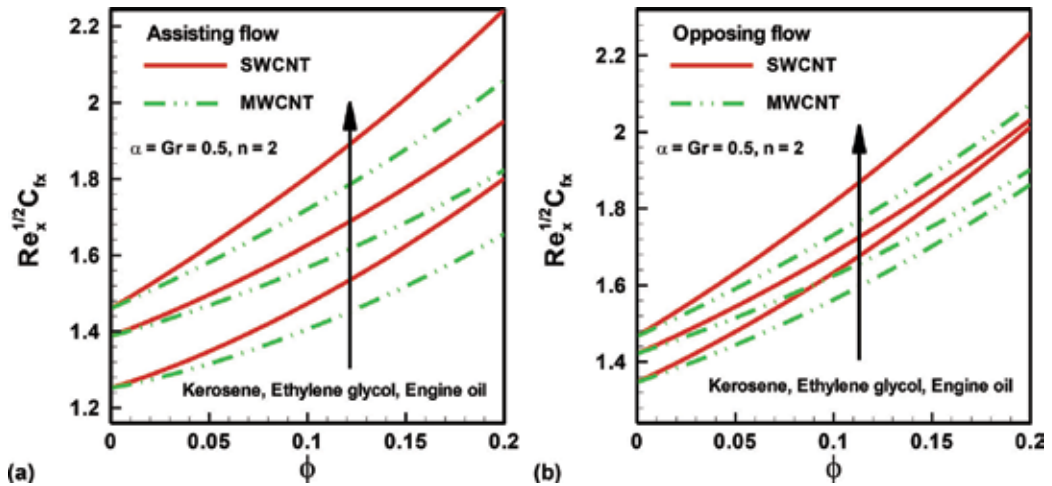


Figure 8. Variation of skin-friction coefficient against nanoparticle volume fraction for different base fluids and different CNTs: (a) assisting flow and (b) opposing flow.

The effects of viscosity parameter (α), power law index (n) and Grashof number (Gr) with various values of nanoparticle volume fraction (ϕ) on temperature profiles are depicted through **Figures 4–6**. For all computations, the nanoparticles and base fluid are considered as MWCNT and kerosene respectively. It is noticed that temperature is maximum at the wall of stretching sheet (i.e. $\eta=0$ or $z=0$). And when we see the temperature surrounding the sheet, it goes down significantly up to $\eta=0.5$ and then after it is constant which is very close to zero. The patterns for temperature profile for assisting and opposing flows are similar. In **Figure 4**, the values of power index and Grashof number are considered as 1 and 0.5 respectively. **Figure 4** depicts that temperature falls with increasing the magnitude of viscosity parameter which is also very clear from the relation $\mu_f(\theta) = e^{-(\alpha\theta)} = 1 - (\alpha\theta) + O(\alpha^2)$ for both type flows (**Figure 4(a and b)**). The values of viscosity parameter and Grashof number are fixed as 0.5 in

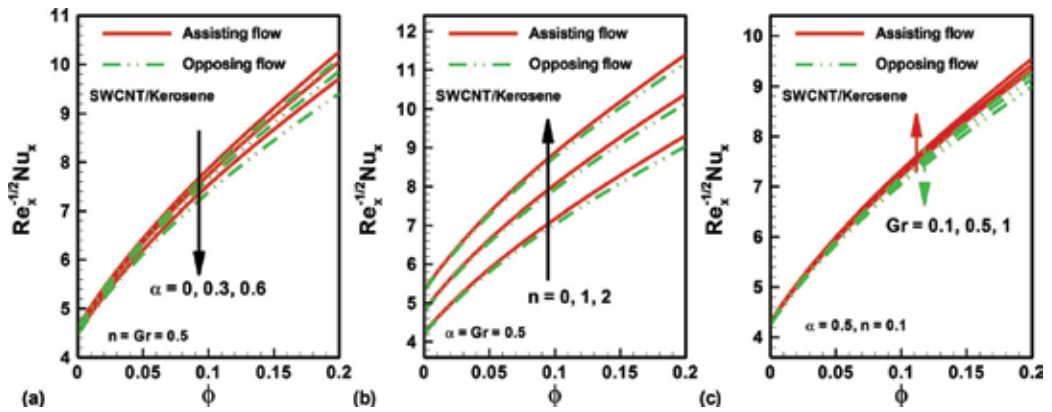


Figure 9. Variation of Nusselt number against nanoparticle volume fraction for (a) viscosity parameter, (b) power law index, and (c) Grashof number.

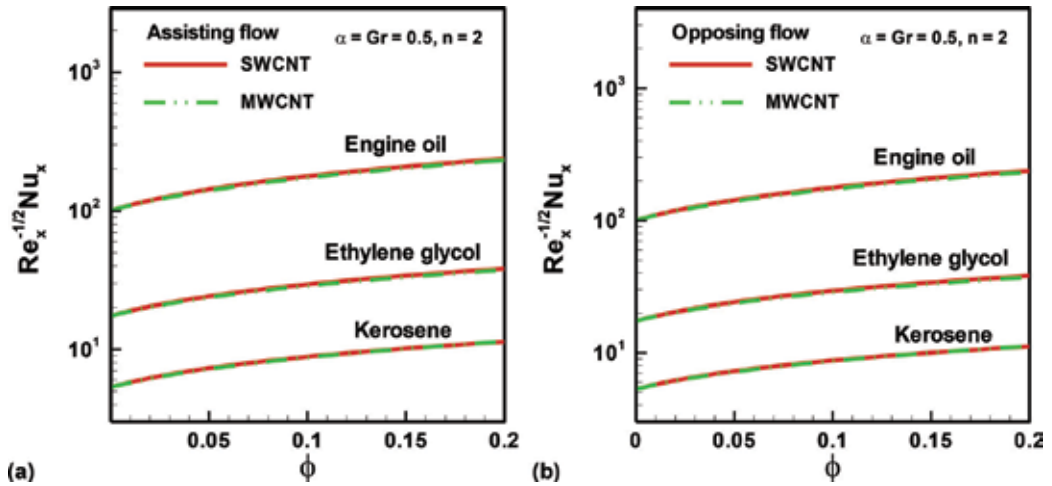


Figure 10. Variation of Nusselt number against nanoparticle volume fraction for different base fluids and also for different CNTs: (a) assisting flow and (b) opposing flow.

Figure 5. It is inferred that the temperature rises with increasing the power law index for different values of $\varphi = 0, 0.1, 0.2$ for both type flows (Figure 5(a and b)). The values of power law index and nanoparticle volume fraction are taken as 0.2 and 0.5 respectively in Figure 6. It is revealed that the temperature slightly increases with increasing the values of Grashof number for assisting flow, see Figure 6a however the temperature decreases with increasing the Grashof number for opposing flow, see Figure 6b. From Figures 4–6, it is also pointed out that the temperature goes up with increasing the nanoparticle volume fraction for all values of viscosity parameter, power law index and Grashof number, and also for both types of flows.

The variations in skin-friction coefficient against nanoparticle volume fraction are computed in Figures 7 and 8 for assisting and opposing flows. It is found that the skin friction coefficient is

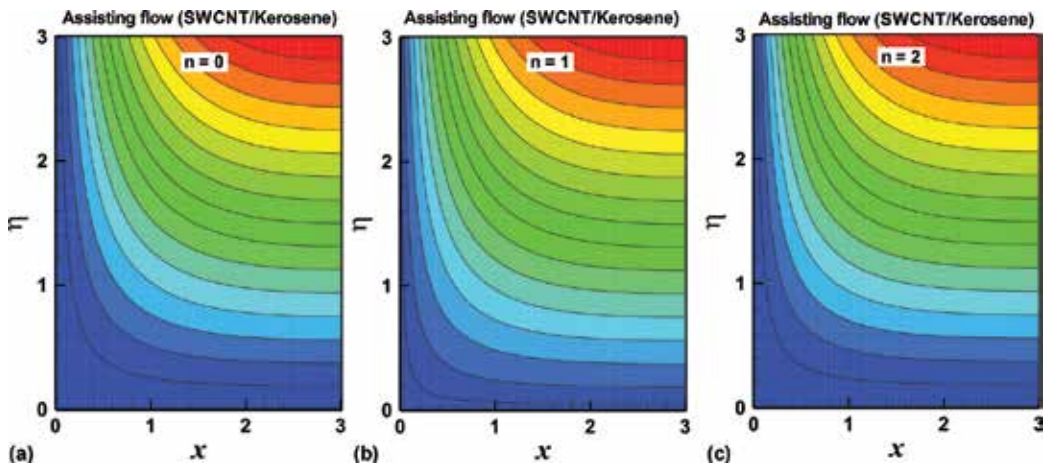


Figure 11. Stream lines for different values of power law index for $\alpha=0.4, Gr=3, \varphi=0.3$.

Physical properties	Base fluid				Nanoparticles	
	Water	Kerosene	Engine oil	Ethylene glycol	MWCNT	SWCNT
ρ (kg/m ³)	997	783	884	1115	1600	2600
c_p (J/kg-K)	4179	2090	1910	2430	796	425
k (W/m-K)	0.613	0.15	0.144	0.253	3000	6600
Pr	6.2	21	6450	2.0363		

Table 1. Thermophysical properties of different base fluid and CNTs.

least at $\varphi=0$ and it starts to enhance nonlinearly with increasing the values of nanoparticle volume fraction. In Figure 7, the nanoparticle and base fluids are considered as SWCNT and Kerosene respectively. The effect of viscosity parameter on skin-friction coefficient at fixed values $n=Gr=0.5$ for both type of flows is presented in Figure 7a. It is observed that the magnitude of skin friction coefficient diminishes with increasing the viscosity parameter for both type of flows. The influence of power law index on skin friction coefficient is examined in Figure 7b at fixed values $\alpha=Gr=0.5$ for both type of flows. It is depicted that skin friction coefficient extends for more values of power law index. Figure 7c is plotted for effects of Grashof number on skin friction coefficient at fixed values $\alpha=n=0.5$ for both type of flows. It is revealed that magnitude of skin friction coefficient diminishes with increasing the magnitude of Grashof number for assisting flow and it increases with increasing the Grashof number for opposing flow.

Figure 8(a and b) are illustrated for the effects of different types of base fluids (kerosene, ethylene glycol and engine oil) and also effects of different CNTs on skin friction coefficient at fixed values $\alpha=Gr=0.5, n=2$ for assisting flow (Figure 8a) and opposing flow (Figure 8b). It is found that the magnitude of skin fraction is more for SWCNT in compare to MWCNT for all type of base fluids

and also for both type of flows. It is further depicted that skin friction is maximum for kerosene and minimum for engine oil for both types of CNTs and also for both types of flows.

Figures 9 and **10** are plotted to see the variations in Nusselt number against nanoparticle volume fraction both types of flows.

It is observed that the magnitude of Nusselt number is minimum at $\varphi=0$. And it significantly increases with increasing the magnitude of nanoparticle volume fraction. The nanoparticle and base fluids are considered as SWCNT and kerosene respectively in **Figure 9**. **Figure 9a** illustrates that Nusselt number diminishes with increasing the magnitude of viscosity parameter for both type of flows at fixed values $n=Gr=0.5$. **Figure 9b** shows that Nusselt number increases with increasing the power index both type of flows at fixed values $\alpha=Gr=0.5$. **Figure 9c** reveals that Nusselt number increases with increasing the Grashof number for assisting flow however it reduces with increasing the Grashof number for opposing flow at fixed values $\alpha=n=0.5$. **Figure 10(a and b)** are plotted to see the effects of Kerosene, Ethylene glycol and Engine oil and also effects of SWCNT and MWCNT on Nusselt number at fixed values $\alpha=Gr=0.5, n=2$ for both types of flows (**Figure 10a and b**). A very minor variation in Nusselt number for SWCNT and MWCNT nanoparticles is noted for assisting and opposing flows. A significant variation in Nusselt number for Kerosene, Ethylene glycol and Engine oil is pointed out where the Nusselt number is largest for Engine oil and smallest for Kerosene for both type of flows.

The stream lines are plotted through the **Figure 11(a-c)** to see the effects of power law index at fixed values $\alpha=0.4, Gr=3, \varphi=0.3$ for assisting flow where SWCNT as nanoparticle and Kerosene as base fluid are considered. It is noted that stream lines go closer to each other with increasing the power law index.

5. Concluding remarks

A numerical investigation for boundary layer flow of CNTs nanofluids with temperature dependent viscosity over a circular stretching sheet is presented. Effects of Three types of base fluids, power law index, viscosity parameter, nanoparticle volume fraction and Grashof number on flow characteristics, and also skin friction coefficient and Nusselt number are discussed appropriately with numerical computations. The concluding remarks of present discussions are précised as:

- The nature of velocity and temperature profiles for assisting and opposing flows are similar for all values of pertinent parameters except Grashof number.
- The velocity boundary layer thickness elaborates with increasing the magnitude of nanoparticle volume fraction, viscosity parameter and power law index for both type of flows.
- The thermal boundary layer thickness expands with reducing the magnitude of viscosity parameter, power law index and nanoparticle volume fraction for both type of flows.
- The magnitude of skin friction coefficient enhances with increasing the nanoparticle volume fraction and power law index and it also enhances with decreasing the viscosity parameter for both type of flows.

- The sequence for skin friction coefficient of different nanoparticles is observed as: $Re_x^{1/2}C_{fx}$ (SWCNT) > $Re_x^{1/2}C_{fx}$ (MWCNT).
- The sequence for skin friction coefficient of different base fluids is also noted as: $Re_x^{1/2}C_{fx}$ (kerosene) > $Re_x^{1/2}C_{fx}$ (ethylene glycol) > $Re_x^{1/2}C_{fx}$ (engine oil).
- The Nusselt number diminishes with increasing the viscosity parameter and it also diminishes with decreasing the power law index and nanoparticle volume fraction for both type of flows.
- The sequence for Nusselt number of kerosene, ethylene glycol and engine oil is noted as: $Re_x^{1/2}Nu_x$ (kerosene) < $Re_x^{1/2}Nu_x$ (ethylene glycol) < $Re_x^{1/2}Nu_x$ (engine oil).

Author details

Noreen Sher Akbar^{1*}, Dharmendra Tripathi² and Zafar Hayat Khan³

*Address all correspondence to: noreensher@yahoo.com

1 DBS&H CEME, National University of Sciences and Technology, Islamabad, Pakistan

2 Department of Mechanical Engineering, Manipal University Jaipur, Rajasthan, India

3 Department of Mathematics, University of Malakand, Dir (Lower), Khyber Pakhtunkhwa, Pakistan

References

- [1] Choi, S. U. S. (1995). Enhancing thermal conductivity of fluids with nanoparticles. ASME, Publications-Fed, 231, 99-106
- [2] Wang XQ, Mujumdar AS. Heat transfer characteristics of nanofluids: A review. International Journal of Thermal Sciences. 2007;**46**(1):1-19
- [3] Das SK, Choi SU, Patel HE. Heat transfer in nanofluids—A review. Heat Transfer Engineering. 2006;**27**(10):3-19
- [4] Mahian O, Kianifar A, Kalogirou SA, Pop I, Wongwises S. A review of the applications of nanofluids in solar energy. International Journal of Heat and Mass Transfer. 2013;**57**(2):582-594
- [5] He Y, Jin Y, Chen H, Ding Y, Cang D, Lu H. Heat transfer and flow behaviour of aqueous suspensions of TiO₂ nanoparticles (nanofluids) flowing upward through a vertical pipe. International Journal of Heat and Mass Transfer. 2007;**50**(11):2272-2281
- [6] Kim SJ, McKrell T, Buongiorno J, Hu LW. Experimental study of flow critical heat flux in alumina-water, zinc-oxide-water, and diamond-water nanofluids. Journal of Heat Transfer. 2009;**131**(4):043204

- [7] Santra AK, Sen S, Chakraborty N. Study of heat transfer due to laminar flow of copper–water nanofluid through two isothermally heated parallel plates. *International Journal of Thermal Sciences*. 2009;**48**(2):391-400
- [8] Ding Y, Alias H, Wen D, Williams RA. Heat transfer of aqueous suspensions of carbon nanotubes (CNT nanofluids). *International Journal of Heat and Mass Transfer*. 2006;**49**(1): 240-250
- [9] Ko GH, Heo K, Lee K, Kim DS, Kim C, Sohn Y, Choi M. An experimental study on the pressure drop of nanofluids containing carbon nanotubes in a horizontal tube. *International Journal of Heat and Mass Transfer*. 2007;**50**(23):4749-4753
- [10] Kumaresan V, Velraj R. Experimental investigation of the thermo-physical properties of water–ethylene glycol mixture based CNT nanofluids. *Thermochimica Acta*. 2012;**545**: 180-186
- [11] Kumaresan V, Khader SMA, Karthikeyan S, Velraj R. Convective heat transfer characteristics of CNT nanofluids in a tubular heat exchanger of various lengths for energy efficient cooling/heating system. *International Journal of Heat and Mass Transfer*. 2013;**60**:413-421
- [12] Walvekar R, Siddiqui MK, Ong S, Ismail AF. Application of CNT nanofluids in a turbulent flow heat exchanger. *Journal of Experimental Nanoscience*. 2016;**11**(1):1-17
- [13] Akbar NS, Kazmi N, Tripathi D, Mir NA. Study of heat transfer on physiological driven movement with CNT nanofluids and variable viscosity. *Computer Methods and Programs in Biomedicine*. 2016;**136**:21-29
- [14] Akbar NS, Abid SA, Tripathi D, Mir NA. Nanostructures study of CNT nanofluids transport with temperature-dependent variable viscosity in a muscular tube. *The European Physical Journal Plus*. 2017;**132**:1-10
- [15] Crane LJ. Flow past a stretching plate. *Zeitschrift für angewandte Mathematik und Physik (ZAMP)*. 1970;**21**(4):645-647
- [16] Khan WA, Pop I. Boundary-layer flow of a nanofluid past a stretching sheet. *International Journal of Heat and Mass Transfer*. 2010;**53**(11):2477-2483
- [17] Makinde OD, Aziz A. Boundary layer flow of a nanofluid past a stretching sheet with a convective boundary condition. *International Journal of Thermal Sciences*. 2011;**50**(7): 1326-1332
- [18] Rana P, Bhargava R. Flow and heat transfer of a nanofluid over a nonlinearly stretching sheet: A numerical study. *Communications in Nonlinear Science and Numerical Simulation*. 2012;**17**(1):212-226
- [19] Xu H, Pop I, You XC. Flow and heat transfer in a nano-liquid film over an unsteady stretching surface. *International Journal of Heat and Mass Transfer*. 2013;**60**:646-652

- [20] Hussain ST, Nadeem S, Haq RU. Model-based analysis of micropolar nanofluid flow over a stretching surface. *The European Physical Journal Plus*. 2014;**129**(8):161
- [21] Uddin MJ, Ferdows M, Bég OA. Group analysis and numerical computation of magneto-convective non-Newtonian nanofluid slip flow from a permeable stretching sheet. *Applied Nanoscience*. 2014;**4**(7):897-910
- [22] Khan WA, Makinde OD, Khan ZH. Non-aligned MHD stagnation point flow of variable viscosity nanofluids past a stretching sheet with radiative heat. *International Journal of Heat and Mass Transfer*. 2016;**96**:525-534
- [23] Hsiao KL. Stagnation electrical MHD nanofluid mixed convection with slip boundary on a stretching sheet. *Applied Thermal Engineering*. 2016;**98**:850-861
- [24] Akbar NS, Tripathi D, Khan ZH, Bég OA. A numerical study of magnetohydrodynamic transport of nanofluids over a vertical stretching sheet with exponential temperature-dependent viscosity and buoyancy effects. *Chemical Physics Letters*. 2016;**661**:20-30
- [25] Bhatti MM, Rashidi MM. Effects of thermo-diffusion and thermal radiation on Williamson nanofluid over a porous shrinking/stretching sheet. *Journal of Molecular Liquids*. 2016;**221**:567-573
- [26] Zeeshan A, Majeed A, Ellahi R. Effect of magnetic dipole on viscous ferro-fluid past a stretching surface with thermal radiation. *Journal of Molecular Liquids*. 2016;**215**:549-554
- [27] Majeed A, Zeeshan A, Ellahi R. Unsteady ferromagnetic liquid flow and heat transfer analysis over a stretching sheet with the effect of dipole and prescribed heat flux. *Journal of Molecular Liquids*. 2016;**223**:528-533
- [28] Hayat T, Muhammad T, Shehzad SA, Alsaedi A. An analytical solution for magnetohydrodynamic Oldroyd-B nanofluid flow induced by a stretching sheet with heat generation/absorption. *International Journal of Thermal Sciences*. 2017;**111**:274-288
- [29] Ellahi R, Hassan M, Zeeshan A. Shape effects of spherical and nonspherical nanoparticles in mixed convection flow over a vertical stretching permeable sheet. *Mechanics of Advanced Materials and Structures*. 2017;**24**(15):1-8
- [30] Nayak MK, Akbar NS, Pandey VS, Khan ZH, Tripathi D. 3D free convective MHD flow of nanofluid over permeable linear stretching sheet with thermal radiation. *Powder Technology*. 2017;**315**:205-215
- [31] Akbar NS, Khan ZH, Nadeem S. The combined effects of slip and convective boundary conditions on stagnation-point flow of CNT suspended nanofluid over a stretching sheet. *Journal of Molecular Liquids*. 2014;**196**:21-25
- [32] Haq RU, Nadeem S, Khan ZH, Noor NFM. Convective heat transfer in MHD slip flow over a stretching surface in the presence of carbon nanotubes. *Physica B: Condensed Matter*. 2015;**457**:40-47

- [33] Hayat T, Hussain Z, Alsaedi A, Asghar S. Carbon nanotubes effects in the stagnation point flow towards a nonlinear stretching sheet with variable thickness. *Advanced Powder Technology*. 2016;**27**(4):1677-1688
- [34] Akbar NS, Khan ZH. Effect of variable thermal conductivity and thermal radiation with CNTS suspended nanofluid over a stretching sheet with convective slip boundary conditions: Numerical study. *Journal of Molecular Liquids*. 2016;**222**:279-286

Numerical Study of Turbulent Flows and Heat Transfer in Coupled Industrial-Scale Tundish of a Continuous Casting Material in Steel Production

Jose Adilson de Castro, Bruno Amaral Pereira,
Roan Sampaio de Souza,
Elizabeth Mendes de Oliveira and
Ivaldo Leão Ferreira

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.75935>

Abstract

This chapter describes the numerical simulations of a coupled industrial scale of the tundish and continuous casting process. The governing equations are presented, and the numerical procedure is discussed in a common framework. The coupled solutions are presented for the transient turbulent flows within the tundish, solidifying zone and extracting regions with the coupling phenomena of heat and mass transfer. The tundish region flow and refractory are calculated using the inlet and outlet boundary conditions in order to estimate the filling phenomena. The transitions and cooling zones for the thin slab continuous casting process are designed to account for the control of the solidified skin in order to avoid breakout. We compared the numerical predictions of the temperatures with industrial monitoring data for a reference case in order to verify the consistence of the model predictions. A parallel version of the numerical code is proposed aiming to improve the computation time keeping numerical accuracy.

Keywords: tundish, continuous casting, turbulent flow, numerical modelling, thin slab, finite volume

1. Introduction

The steel production in an integrated mill requires complex operation units and demands a large amount of energy. In the operation units comprising the transformation of the liquid steel

into slabs, several aspects of the product quality are assured [1–6]. The security and stability of the operations as well as the productivity with lower defects are the main concern and have driven the new development on this step. Of special interest is to fit the dimensions of the slabs suitable for further hot working processes free of internal and superficial defects. The strict control of these steps is the primary effort to high-quality steel slab. A general schematic overview of the continuous casting facilities including the metal transfer steps is presented in **Figure 1**. The initial step is the metal transfer from the ladle to the tundish filling the vessel and establishing the synchronised mass flows. The tundish distributor is used to control the feeding rate of the oscillating mould of the continuous casting step using the submerge tune and flowing valve control. The heat transfer and the flowing phenomena within the oscillating mould are key phenomena to attain the adequate microstructure of the solidified steel and keep the safety of the process with the formation of the solidification skin, which plays the major role on the cooling zones for final solidification of centre of the slab.

The synchronised control of the cooling rate along the mould, bender, speed, and radiation regions is the key for a successful operation of the entire system. In order to improve the process safety, control and productivity comprehensive mathematical models have been developed separately for the tundish and continuous casting processes [1, 4]. Progress in computational simulation has provided tools to help to comprehend the processes. Consequently, several investigations of the parameters which affect the performance under safety operation conditions were driven [3–7]. The tundish operation is carried out in order to assure the compositional and thermal homogeneity of the liquid with low level of impurities and inclusions. The caster machine is designed to promote continuous solidification of liquid metal fed by a tundish through a submerge valve. In the mould region, a strong heat flux is imposed, and a thick solid shell is formed. Water cooling is continuously applied until a secondary region is

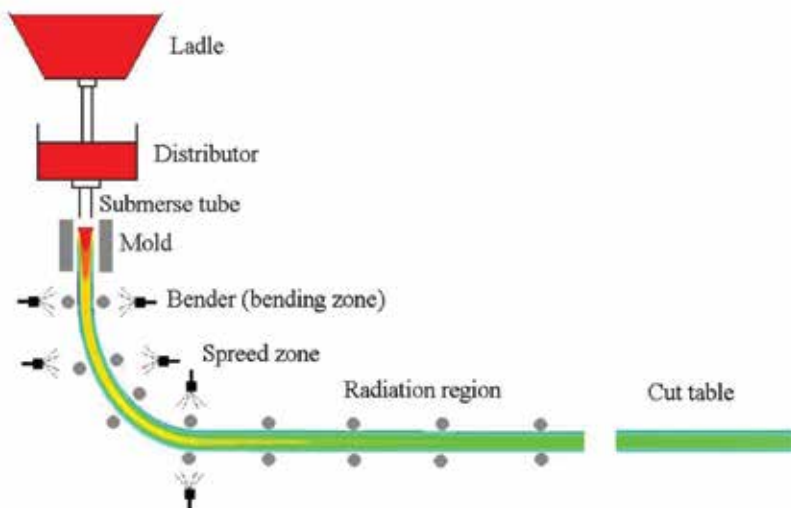


Figure 1. Schematic view of coupling of ladle feeding, tundish and continuous casting process and facilities.

reached, where the cooling is performed only by radiation. At the end of the vein, the slab is cut and discharged on a rolled table. Due to process complexity, which involves heat transfer coupled with phase transformation and fluid flow, the prediction of process parameters and their optimization is usually performed by using empirical procedures. However, the development of efficient numerical techniques and the availability of fast and low-cost computers has burst recently the simulation of real operational conditions [8–13]. To date, it is possible to investigate virtually the manufacturing of several kinds of steels aiming low cost and high material efficiency. Several works in the literature have been reported to analyse the metal behaviour within the oscillating mould of the caster machine due to its importance on the productivity and on the product final quality [4–15]. The oscillating mould is an important component of the caster machine and has strong influence on surface defects and on the temperature distribution inside the mould [2–8]. The heat transfer analysis during solidification is traditionally performed by analytical and numerical methods. Although analytical methods are more elegant, they require a series of assumptions that usually lead to a considerable simplification of the physical phenomena producing unrealistic or limited solutions. Considering numerical methods, four techniques are commonly used: finite differences [7–11], finite elements [12–15], finite volumes [16] and boundary elements [17]. These methods are able to solve the energy, the mass, species and momentum equations. In order to improve scientific calculation performance, continuous changes have been arisen in computational platform paradigm. In the past, the scientific simulation was normally performed in shared memory large computers or in common sequential computers [18]. The fast rate of development in processor technology and the commercial availability of inexpensive powerful personal computers have created a perfect scenario to build up cluster of personal computers as an alternative to the larger and more expensive ones [19]. As consequence of low price, easy maintenance and powerful processors, these so-called Beowulf clusters are becoming popular among scientific computational groups. This architecture offers collective memory to solve scientific complex problems [17, 21]. Although the rise of distributed computer platforms was only an alternative for high-cost supercomputer solutions, they changed profoundly the rule of code development, which now needs to encompass distributed machines [18, 19]. Distributed platforms are suitable for problems in which domain can be split up into small subdomains containing common boundaries. Most CFD codes demand high amount of memory, which is normally available in distributed memory architecture [17, 18]. However, for accuracy and consistency reasons, a parallel implementation needs to interchange information with subdomain boundaries. This synchronisation scheme leads to an increase in the data transfer time due to the existence of a synchronisation elapsed time [17–21]. The communication among computers is carried out by using libraries of Message Passing Interface (MPI) [25]. The library Message Passing Interface (MPI) has largely been used in its freeware version called MPICH [17–21]. In this context, this work newly presents a multidomain parallel numerical model able to simulate the continuous casting of steel. The main objective is to demonstrate the validity of the model and point out the improvement in calculation speedup by developing a code based on multidomain parallel MPI compared to a serial and to a simple MPI parallel code. All the computer codes used in this study are homemade ones, which were developed and tested by the authors.

2. A unified formulation for the tundish and continuous casting processes

The tundish and continuous casting operation units are connected by metal transfer systems to account for the smooth operation and strict control of both. However, a common formulation is possible based on transport phenomena principles. In this section we present a turbulent flow coupled with heat and mass transfer for interconnected processes. The tundish is modelled as a reactor including the metal, slag and inclusions flows, while the refractories and internal protective devices are considered. A multiphase formulation is considered: (a) liquid metal, (b) liquid slag, (c) solidified metal, (d) solidified slag, (e) particle inclusions and (f) refractory:

$$\frac{\partial(\rho u_i)}{\partial t} + \frac{\partial}{\partial x_j}(\rho u_j u_i) = \frac{\partial}{\partial x_j} \left((\mu + \mu_t) \frac{\partial u_i}{\partial x_j} \right) - \frac{\partial}{\partial x_j} (\tau_{ij} + C_{ij} + L_{ij}) - \frac{\partial P}{\partial x_i} - \frac{(1-f_s)}{K_{u_i}} u_i + \rho g_{ui} \quad (1)$$

$$\frac{\partial(\rho T)}{\partial t} + \frac{\partial}{\partial x_j}(\rho u_j T) = \frac{\partial}{\partial x_j} \left(\left(\frac{k}{C_p} + \frac{k_t}{C_p t} \right) \frac{\partial T}{\partial x_j} \right) - \frac{\partial}{\partial x_j} (\theta + C + L) - \frac{\partial}{\partial t} (\rho \Delta H f_s) \quad (2)$$

$$\frac{\partial(\rho)}{\partial t} + \frac{\partial}{\partial x_j}(\rho u_j) = 0 \quad (3)$$

$$\frac{\partial}{\partial t}(\rho C^i) + \frac{\partial(\rho u_j C^i)}{\partial x_j} = \frac{\partial}{\partial x_j} \left(\left(D^i + \frac{\mu_t}{\rho} \right) \frac{\partial(C^i)}{\partial x_j} \right) \quad (4)$$

$$\frac{\partial(\rho k)}{\partial t} + \frac{\partial}{\partial x_j}(\rho u_j k) = \frac{\partial}{\partial x_j} \left(\left(\mu + \frac{\mu_t}{\sigma_k} \right) \frac{\partial k}{\partial x_j} \right) + 2\mu_t S_{ij} S_{ij} - \rho \varepsilon + \beta g_i \frac{\mu_t}{Pr} \frac{\partial T}{\partial x_i} \quad (5)$$

$$\frac{\partial(\rho \varepsilon)}{\partial t} + \frac{\partial}{\partial x_j}(\rho u_j \varepsilon) = \frac{\partial}{\partial x_j} \left(\left(\mu + \frac{\mu_t}{\sigma_\varepsilon} \right) \frac{\partial \varepsilon}{\partial x_j} \right) + C_{1\varepsilon} \frac{\varepsilon}{k} \left(2\mu_t S_{ij} S_{ij} + C_{3\varepsilon} \beta g_i \frac{\mu_t}{Pr} \frac{\partial T}{\partial x_i} \right) - C_{2\varepsilon} \rho \frac{\varepsilon^2}{k} \quad (6)$$

where g_u is the gravity acceleration component at the velocity component direction, depending on the number of species, which can be written as

$$g_u = {}^u g_0 \sum_{C, Mn} \left[\beta_s^i (C_i - C_{l,0}^i) + \beta_T^i (T - T_0) \right] \quad (7)$$

Viscosity was treated as effective viscosity [21] in the following form:

$$\mu_{eff} = \frac{\bar{\sigma}}{3\bar{\varepsilon}} + \mu_t \quad (8)$$

$$\mu_t = (\Delta)^2 \sqrt{S_{ij} S_{ij}} \quad (9)$$

where $\bar{\sigma}$ is the material mean stress and $\bar{\varepsilon}$ is the effective deformation rate presented by Zienkiewicz [22] and Δ is a sub-grid scale for the isotropic turbulence filtering. All the variables used in the formulation are taken as the filtered values. The constants $c_{1\varepsilon} = 1.44$, $c_{2\varepsilon} = 1.92$, $c_{3\varepsilon} = 0.09$ and $\sigma_\varepsilon = 1.30$ are related with turbulent kinetic energy (k) and its dissipation rate (ε):

$$\tau_{ij} = \mu_t \left(\frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i} \right) + \frac{2}{3} k \delta_{ij} \quad (10)$$

$$L_{ij} + C_{ij} = \frac{\Delta k}{12} \left(\frac{\partial u_i}{\partial x_k} + \frac{\partial u_j}{\partial x_k} \right) \quad (11)$$

$$S_{ij} = \frac{1}{2} \left(\frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i} \right) \quad (12)$$

In the solid phase, thermal conductivity was assumed as a function of temperature, according in Holman [23]:

$$k = \psi - 0.01 \gamma T \quad (13)$$

where ψ and γ are constants for a specific metal alloy and the refractories considered for each layer and formed solidified shell. For the liquid phases, a similar relationship is assumed with their specific constants.

The specific heat for the solid and the liquid phases is obtained directly from ThermoCalc calculation using TCFE5 database, while for the refractories, specific relations are used depending on the materials used [24].

The local liquid concentration of each species is given by

$$[C_l^C]_P = \frac{[\rho C]_P^C - [\rho C]_P^{C,old} + [\rho_l g_P^{old} + \beta^C \rho_S (1 - g_P^{old}) k_0^C] [C_l^C]_P^{old}}{\rho_l g_P^{n+1} + \beta^C \rho_S (1 - g_P^{n+1}) k_0^C + (1 - \beta^C) \rho_S k_0^C (g_P^{old} - g_P^{n+1})} \quad (14)$$

$$[C_l^{Mn}]_P = \frac{[\rho C]_P^{Mn} - [\rho C]_P^{Mn,old} + [\rho_l g_P^{old} + \beta^{Mn} \rho_S (1 - g_P^{old}) k_0^{Mn}] [C_l^{Mn}]_P^{old}}{\rho_l g_P^{n+1} + \beta^{Mn} \rho_S (1 - g_P^{n+1}) k_0^{Mn} + (1 - \beta^{Mn}) \rho_S k_0^{Mn} (g_P^{old} - g_P^{n+1})} \quad (15)$$

The segregation parameter β can vary as $0 \leq \beta \leq 1$. Assuming $\beta=1$ means the lever rule, and $\beta=0$, provides Scheil's equation [24].

The heat flux boundary conditions for the continuous caster machine are estimated depending on the region and operational conditions. In the mould and in the foot roll, the cooling water flow is specified at the four faces, internal and external large faces and right and left narrow faces, while at the other zones, heat fluxes were imposed only at two faces, the internal and the external large faces. For both processes the initial conditions are specified by the measured operational conditions or temperature monitoring data.

Heat fluxes on the water-cooled surfaces and on the radiation zones are given by

$$k \frac{\partial T}{\partial x} = h_{eff} (T_{sur} - T_e) + \sigma_r \varepsilon_r (T_{sur}^4 - T_e^4) \quad (16)$$

where h_{eff} is the effective heat transfer coefficient provided by Eq. (17), T_{sur} is the surface temperature, T_e is the environment temperature, σ_r is the Stefan-Boltzmann constant and ε_r is the emissivity.

The heat transfer coefficient in the sprays zones (foot roll, bender and secondary cooling zone) was obtained by the water cooling enthalpy balance, providing

$$h_{eff} = \frac{m_w c_p \Delta T}{A(T_{sur} - T_e)} \quad (17)$$

where m_w is the water flow, c_p is the water-specific heat, A is the cross-sectional area and ΔT is the water temperature difference given as a setup parameter for the cooling system.

The mould region was modelled by using the steel residence time in the mould to calculate the effective heat transfer coefficient. This coefficient regards the effect of thermal resistance due to air gap formation:

$$h_{mold} = 1004.6 \exp(0.02 t_m) \quad (18)$$

where t_m is the steel residence time, calculated by means of the cast velocity setup (V_c) and the mould height (Y) as

$$t_m = \frac{Y}{V_c} \quad (19)$$

The inlet and outlet boundary conditions are specified using mass flow, compositions and temperatures. The wall standard log law is used for modelling the liquid to walls and barrier interfaces in both domains of the tundish and continuous casting vein. **Figure 2** shows the geometry and computational domain of the tundish with the refractories and internal barriers. The numerical mesh used was obtained by using continuous refinement using 20% of the total volume increment of each calculation for a standard operational conditions assuming averaged error less than 1% on the temperature and velocity fields. Same procedure was used to obtain the suitable mesh distribution along the continuous casting vein. The final mesh total volumes in the tundish domain were 201,300 and for the continuous casting vein were 288,000.

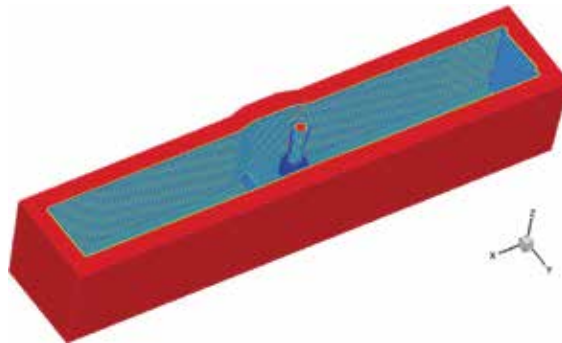


Figure 2. Physical and computational domains including the refractory layers and internal features of the tundish (60 ton).

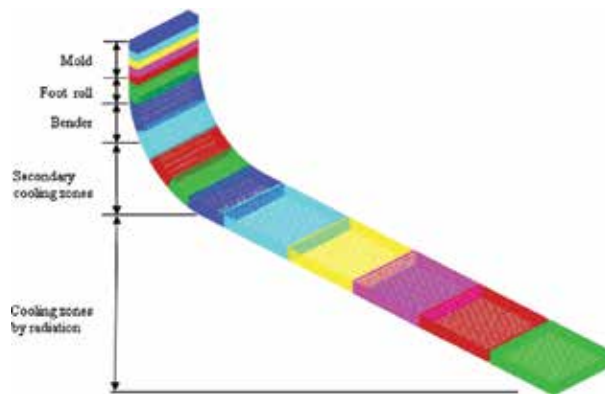


Figure 3. Physical and computational domains indicating the zones of the continuous casting slab and subdomains.

One additional restriction for the continuous casting vein was imposed by using 20% of the total volumes on the oscillating mould region to account for the accuracy of the solution in this region due to the strong gradients developed within the oscillating mould with solidifying shell with strong heat release and solute redistribution. Details of the mesh generated and the subdomains assumed in the simulations are presented in **Figure 3**.

The turbulent quantities at the inlet and outlet are calculated based on the averaged velocities for both processes, as follows:

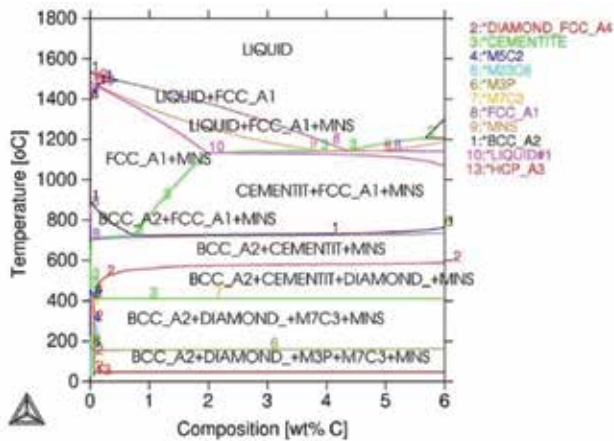
$$U_{av} = \frac{Q}{A} \tag{20}$$

$$k_{av} = 0.01(U_{av})^2 \tag{21}$$

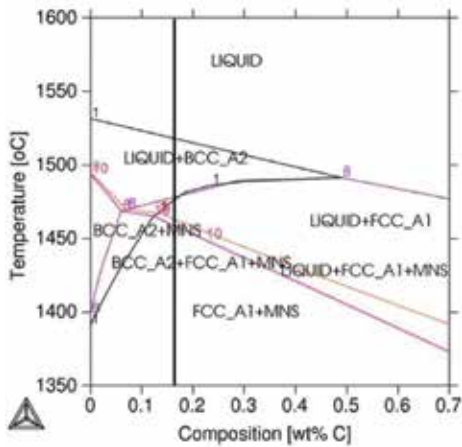
$$\varepsilon_{av} = \frac{2}{D}(k_{av})^{2/3} \tag{22}$$

Eqs. (20)–(22) are applied depending on the geometry of the valves and feeding systems. The averaged values for the temperature and compositions are either set using the measured values or, in the case of transfer system, the values calculated in the previous connected domains.

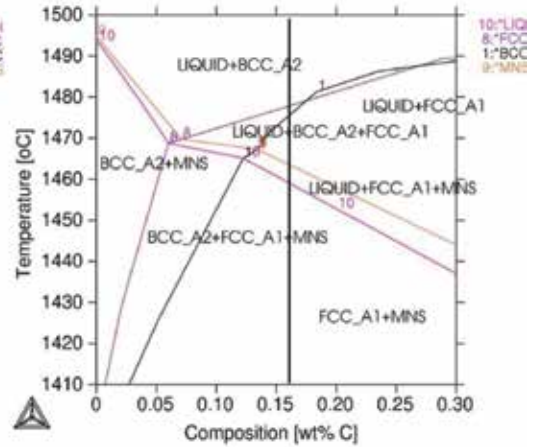
The thermophysical properties of the liquids (steel and slags) and solids formed during the solidification process are determined by using computational thermodynamics database. The solid barriers such as refractories and inhibitors are included by using their tabled thermophysical property data furnished by the suppliers. By using the thermodynamics database, a typical steel is modelled using their pseudo-binary diagrams. **Figure 4** shows the temperatures and phase composition dependency for the whole system and specific regions of the diagram. These data are continuously accessed for local predictions of the thermophysical properties during the transient calculation.



(a)



(b)



(c)

Figure 4. Pseudo-binary phase diagram of Fe-C-Mn-P-S as a function of carbon content (a) phase diagram, (b) and (c) magnification of high temperature range closed to 0, 15 wt% C, which is close to the steel used in this study.

Figure 5 shows the density and heat of phase transformation during temperature evolution. These quantities are accessed to estimate the heat capacity and latent heat released during the solidification and flowing paths. **Figure 6** shows the solid fraction during the solidification path considering the local conditions predicted by computational thermodynamics. With these parameters and physical properties, all the information need for the coupled calculation of the tundish and continuous casting operation are closed.

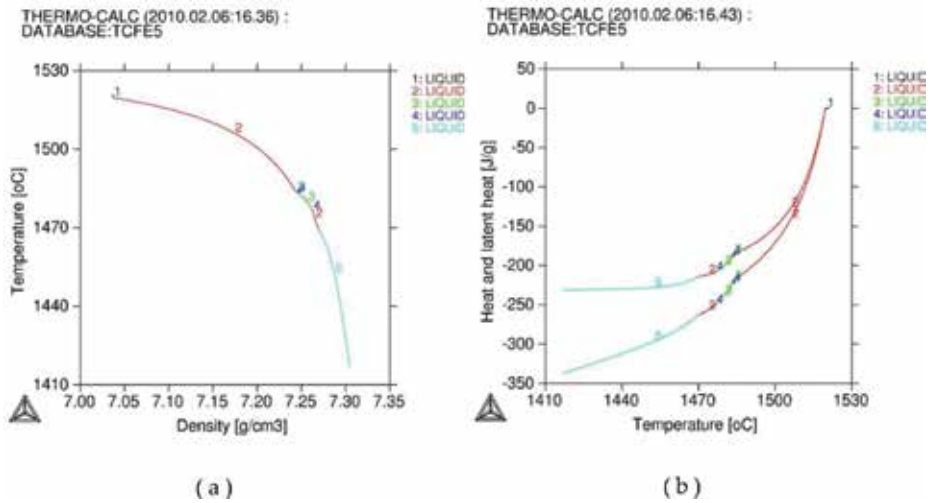


Figure 5. Thermophysical properties inside mushy zone: (a) density and (b) heat and latent heat for the steel used in this computational modelling.

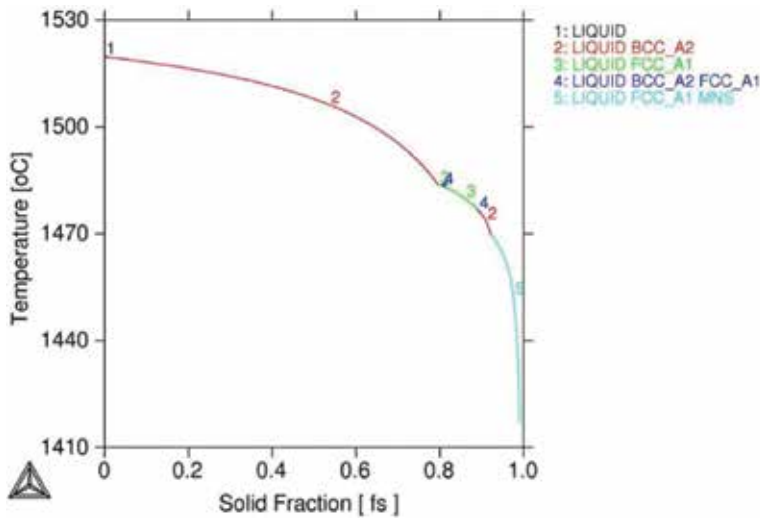


Figure 6. Solid fraction of all solids as a function of temperature during the solidification process within the steel slab used during the calculations accessed from the thermodynamic database.

3. Numerical features

Momentum, mass, energy and species equations were discretized by using the finite volume method (FVM) applied for general coordinate system [25, 26], where the integration is taken

over a typical control volume. The final product of this operation is a set of algebraic equations. Coefficients are obtained by the so-called power law scheme, according to Patankar [26]. The SIMPLE algorithm is used to iteratively determine the velocity components and pressure linked equations. The numerical solution of the set of algebraic equations demands large computational effort. A line-by-line solver based on the tridiagonal matrix algorithm (TDMA) was used to solve the system of algebraic equations. The Alternate Direction Implicit (ADI) iterative procedure was applied within a common solver for all discretized equations. The iterative solution was obtained for each time step in a fully implicit scheme [25, 26]. The convergence criteria were used for all variables admitting a maximum local error less than 1% for all variables simultaneously.

4. Analysis cases

In order to show the capability of the coupled model, a sequence of tundish filling and continuous casting of a Fe-C-Mn thin slab (125 mm) of 1200 mm width is presented. The tundish has 60 ton capacity, and the basic properties of the steel are presented in **Table 1**.

Properties	Units	C-Mn SAE 1018 steel
Slab width	m	1.600
Slab depth	m	0.255
Casting temperature	°C	1.574
Casting speed	m.min ⁻¹	0.810
Cooling water temperature	°C	30
Environment temperature	°C	40
Liquidus temperature	°C	1.519
End of solidification temperature	°C	1.410
Slab material		SAE 1018 steel
Emissivity		0.600
Thermal conductivity in liquid phase	W.m ⁻¹ .K ⁻¹	25.400
Thermal conductivity in solid phase	W.m ⁻¹ .K ⁻¹	29.700
Specific heat in BCC phase at 30–838°C	J.kg ⁻¹ .K ⁻¹	783.400
Specific heat in FCC phase at 838–1416°C	J.kg ⁻¹ .K ⁻¹	647.500
Specific heat in liquid phase 1416–1519°C	J.kg ⁻¹ .K ⁻¹	803.200
Density of solid phase BCC 300–838°C	kg.m ⁻³	7.830
Density of solid phase FCC 838–1416°C	kg.m ⁻³	7.305
Density of liquid at 1522°C	kg.m ⁻³	7.034
Latent heat of solidification, ΔH	J.kg ⁻¹	231.900

Table 1. Basic thermophysical properties of SAE 1018 (Fe-C-Mn) with simulation data.

The initial step of the tundish filling presents strong turbulence features and plays important role on the stable flowing development and security of the whole operation. **Figure 7** shows the flowing pattern ($t=3$ s) for a thin slab operation while the slab extraction is off. As can be observed, the inhibitor apparatus is important to avoid splashing and protect the refractories. **Figure 8** shows the conditions where the stable flow rates are achieved with the liquid level of the tundish nearly constant. The flow pattern indicates that a complex turbulent flow is observed and the liquid flow promotes strong mixing.

In order to assure the coupled model formulation and the simultaneous solution in the parallel platform simulation, a confrontation with measured temperature profile measured in the industrial machine was performed. **Figure 9** showed a comparison of the model predictions for the serial, parallel and the pyrometer measurement at the industrial machine.

As can be observed, a close agreement with the industrial operation measurements for the temperature is reached. The measurements and calculations were compared for the stable casting operation, and the measured values were obtained using infrared pyrometer, and the plotted values are the average of five runs with intervals of 5 min.

As can be observed also, the values obtained with the serial and parallel versions are virtually the same. A complete view of the solid portion of the continuous casting domain is shown in **Figure 10** for stable flowing state. A thin skin formed in the oscillating mould region and continuous growing along the bending and cooling zones is observed. A recalescence and final cooling regions are observed. These regions are critical for the process due to the possibility of crack and defect susceptibility depending on the cooling rates and inclusions dragged and formed during the casting development [22–24].

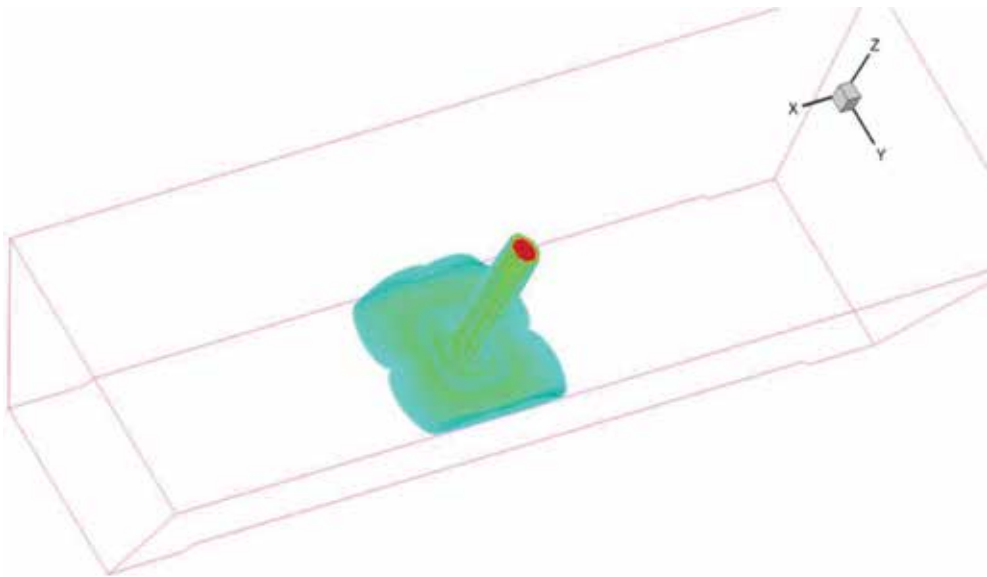


Figure 7. Fluid flow pattern during initial stage of the tundish filling period.

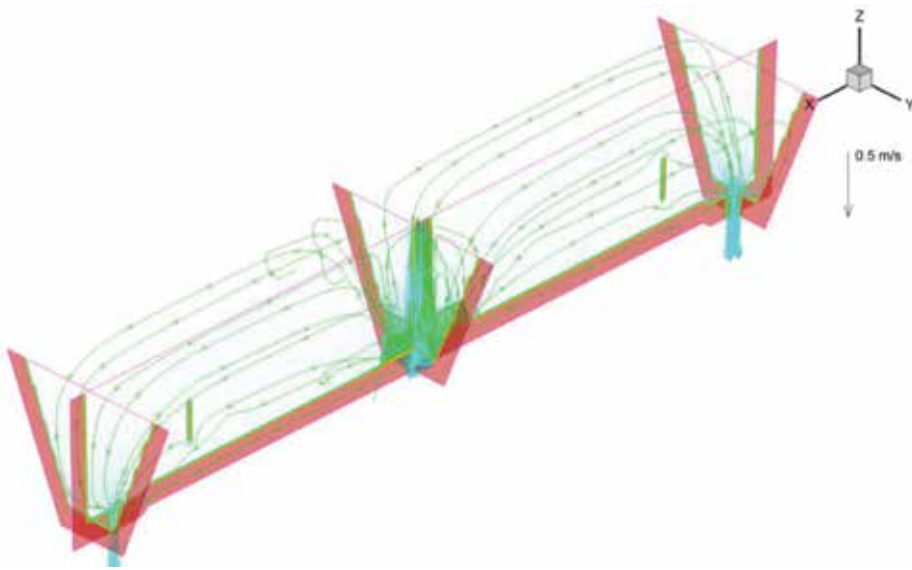


Figure 8. Fluid flow pattern obtained with the solutions of the model equations with parallel code version.

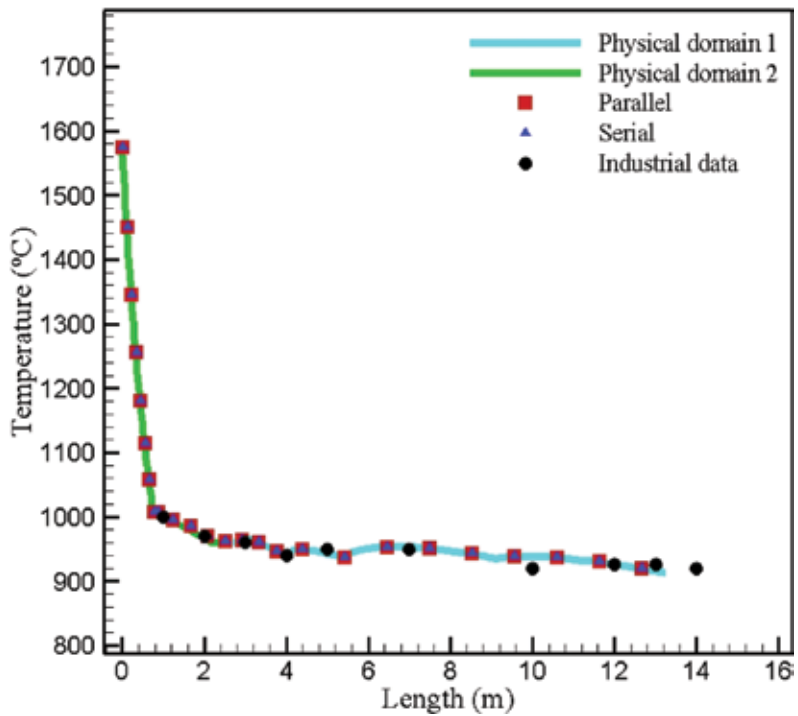


Figure 9. Model validation with industrial data acquisition along the steel slab and comparison with the solutions obtained with the serial and parallel code versions.

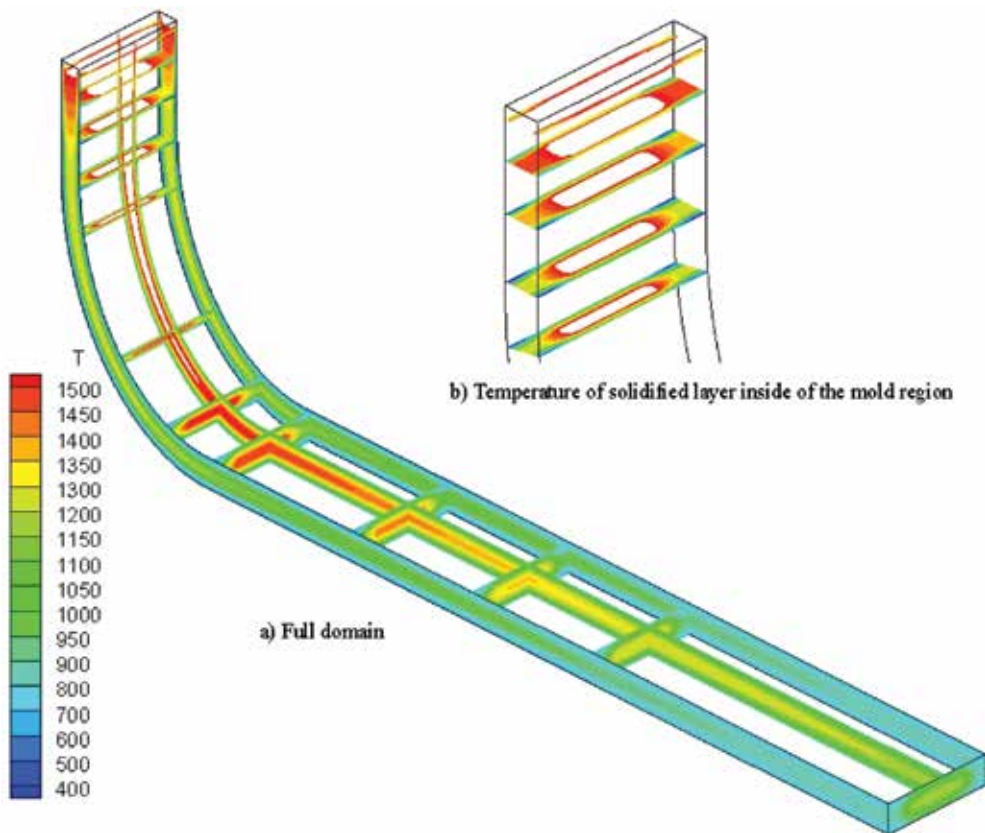


Figure 10. Full domain temperature distribution pattern and the solid skin formed during the continuous casting process of steel slab.

5. Conclusions

A unified formulation for the liquid metal flows and heat transfer within the tundish and continuous casting was presented and applied for actual industrial practices. New operational conditions for the metal flows aiming to allow the inclusion flotation and slag capture are suggested. The prediction of actual continuous casting practice is compared with industrial data. Thus, the simulation platform can be used for designing new thin slab continuous casting process, which could decrease the subsequent steps of hot rolling for thickness reduction.

Acknowledgements

Authors acknowledge financial support provided by FAPERJ (Fundação de Amparo à Pesquisa do Estado do Rio de Janeiro), CAPES (Coordenação de Aperfeiçoamento de Pessoal de Ensino superior) and CNPq (Conselho Nacional de Pesquisa)

Author details

Jose Adilson de Castro^{1*}, Bruno Amaral Pereira¹, Roan Sampaio de Souza¹,
Elizabeth Mendes de Oliveira² and Ivaldo Leão Ferreira³

*Address all correspondence to: joseadilsoncastro@id.uff.br

1 Federal Fluminense University (UFF/RJ), Brazil

2 Department of Metallurgical Engineering, Federal Center for Technological Education (CEFET/RJ), Angra dos Reis, Brazil

3 Department of Mechanical Engineering, Federal University of Pará (UFPA/Pa), Belém, Brazil

References

- [1] Lan XK, Khodadadi JM. Fluid flow, heat transfer and solidification in mold of continuous caster during ladle change. *International Journal of Heat and Mass Transfer*. 2001;**44**:953-965. DOI: 10.1016/S0017-9310(00)00145-9
- [2] Wang EG, He JC. Finite element numerical simulation on thermo-mechanical behavior of steel billet in continuous casting mold. *Science and Technology of Advanced Materials*. 2001;**2**:257-263. DOI: 200.130.19.152
- [3] Luo X, Xie Q, Wang Y, Yang C. Estimation of heat transfer coefficients in continuous casting under large disturbance by Gaussian kernel particle swarm optimization method. *Cuie International Journal of Heat and Mass Transfer*. 2017;**111**:1087-1097. DOI: 10.1016/j.ijheatmasstransfer.2017.03.105
- [4] Pardeshi R, Basak S, Singh AK, Basu B, Mahashabde V, Roy SK, Kumar S. Mathematical modeling of the Tundish of a single-strand slab caster. *ISIJ International*. 2004;**44**:1534-1540. DOI: 10.2355/isijinternational.44.1534
- [5] Ha MY, Lee HG, Seong SH. Numerical simulation of three-dimensional flow, heat transfer and solidification of steel in continuous casting mold with electromagnetic brake. *Journal of Materials Processing Technology*. 2003;**133**:322-339. DOI: 10.1016/S0924-0136(02)01009-9
- [6] Janik M, Dya H. Modelling of three-dimensional temperature field inside the mould during continuous casting of steel. *Journal of Materials Processing Technology*. 2004;**157**:177-182. DOI: 10.1016/j.jmatprotec.2004.09.026
- [7] Peng X, Zhou J, Qin Y. Improvement of temperature distribution in continuous casting moulds through the rearrangement of the cooling water slots. *Journal of Materials Processing Technology*. 2005;**167**:508-514. DOI: 10.1016/j.jmatprotec.2005.05.023
- [8] Choudhary SK, Mazumdar D, Ghosh A. Mathematical modeling of heat transfer phenomena in continuous casting of steel. *ISIJ International*. 1993;**33**:764-774. DOI: 0915-1559

- [9] Zhou L, Wang W, Xu C, Chen Z. An investigation of the mold-flux performance for the casting of Cr12MoV steel using a mold simulator technique. *Metallurgical and Materials Transactions*. 2017;**48**:2017-2026. DOI: 10.1007/s11663-017-0990-0
- [10] Maurya A, Jha - Pradeep K. Influence of electromagnetic stirrer position on fluid flow and solidification in continuous casting mold. *Applied Mathematical Modelling*. 2017;**48**:736-748. DOI: 10.1016/j.apm.2017.02.029
- [11] Spinelli JE, Tosetti JP, Santos CA, Spim JA, Garcia A. Microstructure and solidification thermal parameters in thin strip continuous casting of stainless steel. *Journal of Materials Processing Technology*. 2004;**150**:255-262. DOI: 10.1016/j.jmatprotec.2004.02.040
- [12] Thomas BG, Mika LJ, Najjar FM. Simulation of fluid flow inside a continuous slab-casting machine. *Metallurgical Transactions B*. 1990;**21B**:387-400. DOI: 0360-2141
- [13] Chan YW. Finite element simulation of heat flow in continuous casting. *Advanced Engineering Software*. 1989;**11**:128-135. DOI: 10.1016/0141-1195(89)90042-9
- [14] Brian GT, Fady MN. Finite element modeling of turbulent fluid flow and heat transfer in continuous casting. *Applied Mathematical Modelling*. 1991;**15**:226-243. DOI: 10.1016/0307-904X(91)90001-6
- [15] Sheng-Long J, Zheng Z, Liu M. A multi-stage dynamic soft scheduling algorithm for the uncertain steelmaking-continuous casting scheduling problem. *Applied Soft Computing*. 2017;**60**:722-736. DOI: 10.1016/j.asoc.2017.07.016
- [16] Husepe AE, Cardona A, Fachinotti V. Thermomechanical model of continuous casting process. *Computer Methods in Applied Mechanics Engineering*. 2000;**182**:439-455. DOI: 10.2478/v10172-012-0034-3
- [17] Vasta VN, Hammond DP. Viscous flow computations for complex geometries on parallel computers. *Advances in Engineering Software*. 1998;**29**:337-343. DOI: 10.1016/S0965-9978(97)00076-8
- [18] Sikora J, Ramakrishnan S, Leblanc L. Conversion of a single process CFD code to distributed and massively parallel processing. *Advanced Engineering Software*. 1998;**29**:331-336. DOI: 10.1016/S0965-9978(98)00007-6
- [19] Lepper J, Schnell U, Hein KRG. Parallelization of a simulation code for reactive flows on the intel paragon. *Computers & Mathematics with Applications*. 1998;**35**:101-109. DOI: 10.1016/S0898-1221(98)00037-6
- [20] Nesterov O. A simple parallelization technique with MPI for ocean circulation models. *Journal of Parallel and Distributed Computing*. 2010;**70**:35-44. DOI: 10.1016/j.jpdc.2009.09.005
- [21] Moreira LP, Castro JA. Modeling the hot rolling process using a finite volume approach. *Transactions on Engineering Sciences*. 2008;**59**:419-430
- [22] Zienkiewicz OC. Flow of solids during forming and extrusion: Some aspects of numerical solutions. *International Journal of Solids and Structures*. 1978;**14**:15-38. DOI: 10.1016/0020-7683(78)90062-8

- [23] Holman JP. Heat Transfer. 7th ed. New York, NY: Mc Graw-Hill Book Company; 1990
- [24] Ferreira IL, Voller VR, Nestler B, Garcia A. Two-dimensional numerical model for the analysis of macrosegregation during solidification. *Computational Materials Science*. 2009;**46**:358-366. DOI: 10.1016/j.commatsci.2009.03.020
- [25] Melaaen MC. Calculation of fluid flows with staggered and nonstaggered curvilinear nonorthogonal grids—The theory. *Numerical Heat Transfer-Part B*. 1992;**21**:1-19
- [26] Patankar VS. *Numerical Heat and Fluid Flow*. 2nd ed. Washington, DC: Hemisphere Publishing Corp; 1980. 197 p. DOI: 007084740-5

Modeling of the Temperature Field in the Magnetic Hyperthermia

Iordana Astefanoaei and Alexandru Stancu

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.71364>

Abstract

The numerical and/or analytical modeling of the temperature field developed by the magnetic systems in the external alternating magnetic fields is essential in the Magnetic Hyperthermia. Optimization of the all parameters involved in the burning process of the malignant tissues can be realized more efficiently using a mathematical model. The analytical models can be used for the validation of any numerical complex models of the heating processes. This work focuses on the parameters which influences the therapeutic temperature field developed by the magnetic systems within the malignant tissues when the magnetic field is applied. An analytical model was developed to predict and control the bioheat transport within a malignant tissue. This model was compared with a numerical model which was developed in the same conditions of the thermal analysis. Infusion of a diluted suspension of magnetic nanoparticles (MNP) into liver tissue was modeled using the Darcy's equation. The MNP concentration and the temperature field were computed for different parameters as: (i) ferrofluid infusion rates, (ii) particle zeta potential and (iii) magnetic field parameters. The convection-diffusion-deposition of the particles within tissues was considered in this analysis. This study indicates the essential role of these parameters to predict accurately the hyperthermic temperature field. The model presented in this paper predicts the optimum MNP dosage and the temperature at every point within the malignant tissue.

Keywords: finite element method (FEM), Magnetic Hyperthermia, modeling of the temperature field, magnetite nanoparticles

1. Introduction

The Magnetic Hyperthermia is one of the most promising therapies in the cancer treatment [1]. The malignant tissues are destroyed when their temperature reach a therapeutic

hyperthermic temperature in the range 40–46°C [2, 3]. The main problem of this technique is to understand and to control as precisely as possible the temperature field developed by the magnetic systems injected within malignant tissues when the external alternating magnetic fields are applied.

Some experimental data realized in a tissue equivalent (agarose gel) evidences the particle diffusion within the tissue after their injection [4, 5]. The diffusion-convection and deposition of the particles have a strong influence on the radial particle distribution within the tissue volume and as a consequence on the temperature field in tissues [6].

In this paper, the temperature within a malignant tissue surrounded by a healthy tissue was studied considering the radial magnetic nanoparticles (MNP) concentration as an effect of the ferrofluid injection at the center of tumor. The MNP with different sizes having a lognormal particle size distribution were considered. The temperature developed by the magnetic systems in the external time-dependent magnetic field was analyzed for different values of the parameters. During the injection process of the particles within the tissues, their convection and deposition influences strongly the concentration of the particles. An analytical model was developed to predict the temperature field for different important parameters as: (i) ferrofluid infusion rates, (ii) particle zeta potential and (iii) optimal particle dosages. The results were compared with a numerical model in Comsol Multiphysics and Matlab. The values of temperatures computed using the analytical and numerical models in the same conditions were in very good agreement.

2. The analytical temperature model

Modeling the heat transport within the malignant tissues is an important element in the Magnetic Hyperthermia. The development of an analytical model for a complex system or process is usually an important breakthrough with a strong impact in the field.

In this case we are presenting an analytical model which will be a powerful analysis tool for hyperthermia treatment planning due to its ability to predict accurately the temperature field within the malignant tissues. The model will provide a tool to study the main parameters which influence significantly the temperature field and to give a tool to optimize them. For each patient, an individual therapy planning is required in correlation with the tumor location, geometry, shape and size. The elaboration of a patient model by segmentation of images from computerized tomography or magnetic resonance imaging scans is the first step—and the most important one—of hyperthermia treatment planning. This segmentation is used to compute the temperature field in the tumor located in a patient organ.

In our simulations a spherical concentric configuration composed of a malignant tissue surrounded by a healthy tissue was considered (**Figure 1**). The malignant tissue has a radius R_1 and the healthy region has the shell thickness $R_2 - R_1$. At the center O of this geometric structure,

a ferrofluid volume V_f was injected with the volumetric flow rate (ferrofluid infusion rate) Q_v ($\mu\text{l}/\text{min}$) using a needle of a syringe with the radius r_o . The ferrofluid flow within this geometry starts from the center O where is localized the injection site (IS).

In this analysis, MNP with different sizes were considered. A lognormal distribution was defined by the following distribution function [7]:

$$g[R] = \frac{1}{\sigma R \sqrt{2\pi}} \exp\left[-\frac{[\ln(R/R_0)]^2}{2\sigma^2}\right] \text{ and } \int_0^\infty g[R] dR = 1$$

R is the particle radius, $\ln[R_0]$ is the median and σ is the standard deviation of R.

2.1. The radial distribution of the MNP concentration

The ferrofluid (composed by small magnetite particles and water) was considered an incompressible diluted colloidal fluid with the small concentration ($c \leq 5\%$ by volume). The presence of the small magnetite particles does not significantly affect the transport properties of the fluid [8]. The velocity of the magnetic particles within tissues was computed as a solution of the continuity equation in the spherical coordinates:

$$\nabla \cdot \vec{v} = 0 \tag{1}$$

The radial velocity of the particles—the component of the velocity vector: $\vec{v}(v, 0, 0)$ is given by:

$$v_r = \frac{B}{r^2} \tag{1.1}$$

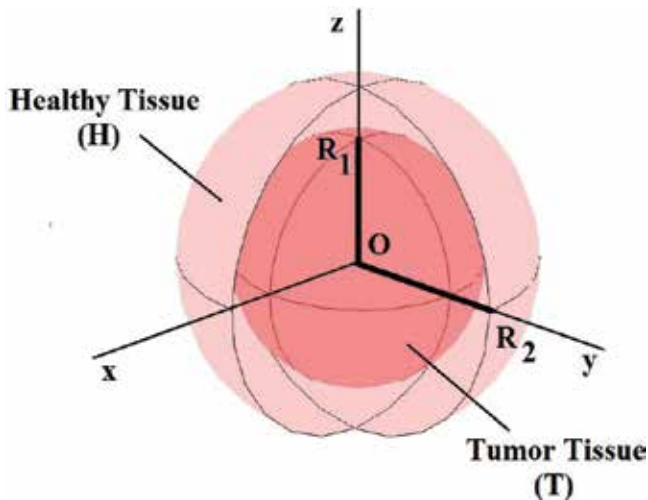


Figure 1. The geometric configuration of a malignant and healthy tissue structure.

where the variable r defines the radial distance from the IS localized at the center of this geometry. The constant $B = \frac{Q_v}{\pi}$ was computed considering the ferrofluid velocity at the tip of the needle as:

$$U = \frac{Q_v}{S_{\text{needle}}} = \frac{B}{r_0^2} \quad (1.2)$$

$S_{\text{needle}} = \pi r_0^2$ is the needle cross sectional area. The radial velocity of the particles depends on the volumetric flow rate Q_v ($\mu\text{l}/\text{min}$) and the radial variable r :

$$v_r = \frac{Q_v}{\pi r^2} \quad (2)$$

Local velocity of the particles (2) depends on the pressure gradient developed within geometry as result of the ferrofluid injection process. The ferrofluid flow through tissues was modeled using the Darcy's equation [8]:

$$\nabla P_i = -\frac{\varepsilon_i \mu}{K_i} \vec{v} \quad (i = 1, 2) \quad (3)$$

In this analysis, the index $i = 1$ defines the tumor (malignant tissue) and the index $i = 2$ defines the healthy tissue. The expression of the pressure in the malignant tissue is $P_1(r)$:

$$P_1(r) = \frac{\varepsilon_1 \mu}{K_1} \left[\frac{1}{r} - \frac{1}{R_1} \right] \frac{Q_v}{\pi} + \frac{\varepsilon_2 \mu}{K_2} \left[\frac{1}{R_1} - \frac{1}{R_2} \right] \frac{Q_v}{\pi}, \quad r_0 \leq r \leq R_1. \quad (4.1)$$

and $P_2(r)$ in the healthy region:

$$P_2(r) = \frac{\varepsilon_2 \mu}{K_2} \left[\frac{1}{r} - \frac{1}{R_2} \right] \frac{Q_v}{\pi} \quad R_1 \leq r \leq R_2 \quad (4.2)$$

was computed solving the Darcy's equation (3) for the concentric tissues. On the external border of the geometry, the pressure is zero, $P_2(r=R_2)=0$. The pressure developed in this geometric configuration depends significantly on the parameter Q_v , ferrofluid viscosity μ and tissues characteristics (porosity ε_i and permeability K_i).

The mass concentrations of the particles, $C_i = C_i(r)$ ($i = 1, 2$) (expressed in mg/cm^3) are the solutions of the modified convection-diffusion equation [8, 9]:

$$\frac{\partial C_i}{\partial t} + \nabla \cdot (\vec{u}_i C_i) = \nabla \cdot (D_i^* \nabla C_i) - k_f^i C_i \quad (5)$$

where k_f^i represent the values of the deposition rate coefficients of the particles within the malignant and healthy tissues which were computed using the relations (A1.3)–(A1.5) from **Appendix 1**.

The expression $\nabla \cdot (\vec{u}_i C_i)$ describes the particle convection, the term $\nabla \cdot (D_i^* \nabla C_i)$ - the particle diffusion and $k_f^i C_i$ - the mean volumetric deposition rate of the particle on the solid phase.

The deposition processes play an important role in the spatial distribution of the particles. The temperature field within geometry is strongly dependent on the volumetric deposition rate of the particles on the solid phase $k_f^i C_i$. The coefficients k_f^i

$$k_f^i = \frac{3(1 - \varepsilon_i)}{2 \varepsilon_i d_c} \eta_s^i v_r \quad (i = 1, 2)$$

were computed considering the superposition of the effects developed by the hydrodynamic forces, van der Waals interactions, gravity effect and the repulsive electrostatic double layer forces. The particle deposition on the cellular structure of tissues depends on the particle diameter D , the radial velocity v_r of the particles, tissues characteristics (porosity ε_i and permeability K_i) and particle zeta potential ζ_p (**Appendix 1**) [10–16]. At equilibrium, Eq. (5) have the following solutions (**Appendix 1** - relations (A1.15)):

$$C_i(r) = \left(\frac{e^{-\frac{A_i}{2r}}}{\sqrt{r}} \right) \left\{ (\text{const1})_i \text{Bessel I} \left[\sqrt{1 - 4 m_i}, \frac{A_i}{2r} \right] + (\text{const2})_i \text{Bessel K} \left[\sqrt{1 - 4 m_i}, \frac{A_i}{2r} \right] \right\} \quad (6)$$

where the expressions m_i and A_i are:

$$m_i = -\frac{3(1 - \varepsilon_i)}{2 \varepsilon_i d_c} \eta_s^i \frac{Q_v}{\pi D_i^*} \text{ and } A_i = \frac{Q_v}{\pi D_i^*}.$$

The collector efficiency η_s^i as result of the electrostatic repulsive forces is given by the relation (A1.2) from **Appendix 1**.

Boundaries conditions: The constants, $(\text{const1})_i$ and $(\text{const2})_i$, were computed using the following boundary conditions:

- (i) $C_2 = 0$ on the external boundary of the geometry ($r = R_2$);
- (ii) Neumann boundary condition at the all inner interfaces:

$$C_1(r = R_1) = C_2(r = R_1)$$

$$D_1^* \frac{\partial C_1}{\partial r} \Big|_{r=R_1} = D_2^* \frac{\partial C_2}{\partial r} \Big|_{r=R_1}$$

- (iii) at the injection site (IS), at the top of the needle ($r = r_0$) the concentration has the particular expression $C_1 = C_{\max}$.

The radial distribution of the MNP concentrations computed as solutions of Eq. (5) depends mainly on: (i) the ferrofluid infusion rate Q_v ; (ii) tissue and particle characteristics as porosity, permeability and particle size (**Table 1**). The convection, diffusion and deposition of the

Hamaker constant A (J)	$3 \cdot 10^{-21} - 4 \cdot 10^{-20}$
Particle radius R (nm)	5–30
Collector diameter d_c (mm)	0.05–0.50
Tissue porosity: $\varepsilon_i (i=1, 2)$ ε_1 —malignant tissue porosity; ε_2 —healthy tissue porosity	$\varepsilon_1 = 0.1-0.8$; $\varepsilon_2 = 0.2$
Permeability $K_i (i=1, 2)$ K_1 —malignant tissue permeability; K_2 —healthy tissue permeability	$K_1 = 10^{-14}$; $K_2 = 5 \cdot 10^{-13}$;
Water mass density ρ_w (kg/m ³)	1000
Boltzmann coefficient k_b (J/K)	$1.38 \cdot 10^{-23}$
Absolute fluid viscosity μ (kg/(s m))	0.001
ζ_p —particle zeta potential (mV)	-10 to -50
ζ_c —collector zeta potential (mV)	-20

Table 1. Parameters values used in simulations [9, 10].

particles influence strongly the spatial distribution of the particles after their injection within tissues. At equilibrium, the temperature field depends significantly on the spatial distribution of the particles. In this analysis, the volume fractions of the particles were considered:

$$\Phi_i(r) = \frac{C_i(r)}{\rho_{MNP}} \quad (7)$$

ρ_{MNP} is the mass density of the magnetic particles. At the injection site - at the center of the geometry - the maximum value of the volume fraction was $\Phi_{\max} = \frac{C_{\max}}{\rho_{MNP}}$.

2.2. The temperature field

The temperature field within malignant and healthy tissues is described by the solutions: $T_i = T_i(r, t)$, ($i=1, 2$) of the bioheat transfer equation (Pennes equation) in the living tissues [6]:

$$\rho_i c_i \frac{\partial T_i}{\partial t} = \nabla [k_i \nabla T_i] + \rho_b w_b c_b [T_{art} - T_i] + Q_{met}^i + Q_i(r), \quad (8)$$

with the following thermal characteristics (**Table 2**): ρ_i —the mass density, c_i —specific heat capacity, k_i —the thermal conductivity, ρ_b —mass density of the blood, c_b —specific heat capacity of the blood, T_{art} —blood temperature, Q_{met}^i —metabolic heat production and $Q_i(r)$ (W/m³)—power density (volumetric heating rate) dissipated by the magnetic particles within geometric configuration when the magnetic field is applied.

As a result of the spatial distribution of the particles, the total volumetric heating rate $Q_i(r)$ depends on the radial dependent volume fractions of the particles $\Phi_i(r)$. When MNP with different sizes are injected within tissues, the volumetric heating rate of the ferrofluid is [7]:

Thermal and magnetic characteristics	Magnetite	Tumor tissue	Healthy tissue
Mass density (kg/m ³)	5180	1160	1060
Specific heat capacity (J/kg K)	670	3600	3600
Thermal conductivity (W/mK)	40	0.4692	0.512
Magnetization (kA/m)	446	–	–
Anisotropy constant K (kJ/m ³)	9	–	–
Frequency f (kHz)	100–650	–	–
Magnetic field amplitude H (kA/m)	0–15	–	–

Table 2. The thermal and magnetic characteristics [9, 10].

$$Q_i(r) = \Phi_i(r) \bar{P} \tag{9}$$

The size distribution of particles influences strongly the heating rate. For a magnetic system which contains the particles with different sizes, the volumetric heating rate is given by

$$\bar{P} = \int_0^\infty P[R]g[R]dR,$$

where the volumetric heating rate released by one particle $P[R]$ and susceptibility $\chi'' [R]$ depend strongly on the particle radius:

$$P[R] = \mu_0 f \pi \chi'' [R] H_0^2 \text{ and } \chi'' [R] = \frac{\mu_0 M_s^2 V[R]}{3 k_B T} \frac{2\pi f \tau[R]}{1 + (2\pi f \tau[R])^2} \tag{10}$$

H_0 is the intensity of the magnetic field, f —frequency of the magnetic field, $\mu_0 = 4\pi \cdot 10^{-7}$ H/m the permeability, M_s is the saturation magnetization and $k_B = 1.38 \cdot 10^{-23}$ J/K is Boltzmann constant. The effective relaxation time contains the Brown relaxation time $\tau_B [R]$ and the Néel relaxation time $\tau_N [R]$ as function of the particle radius:

$$\tau [R] = \frac{\tau_N [R] \tau_B [R]}{\tau_N [R] + \tau_B [R]}$$

$$\tau_B [R] = \frac{3 \eta V_H [R]}{k_B T}; \tau_N [R] = \tau_0 \frac{\sqrt{\pi}}{2} \frac{\exp\left[\frac{K V [R]}{k_B T}\right]}{\sqrt{\frac{K V [R]}{k_B T}}}.$$

τ_0 is the average relaxation time, η carrier liquid viscosity, $V_H [R] = V [R] \left(1 + \frac{\delta}{R}\right)^3$ is the hydrodynamic volume of the particles, δ is the surfactant layer thickness and K is the anisotropy constant of the magnetic particles. Using the spherical symmetry, Eq. (8) can be written as [17]:

$$\frac{k_i}{r^2} \frac{\partial}{\partial r} \left[r^2 \frac{\partial T_i}{\partial r} \right] + M_i [T] = \rho_i c_i \frac{\partial T_i}{\partial t} \tag{11}$$

with $M_i[T] = \omega_b^i \rho_b c_b [T_e^i(r) - T_i]$. For $i = 1, 2$, the expressions $T_e^i(r)$ are:

$$T_e^1(r) = T_B^1 + \frac{Q_1(r)}{\omega_b^1 c_b \rho_b} \text{ and } T_e^2(r) = T_B^2 + \frac{Q_2(r)}{\omega_b^2 c_b \rho_b}$$

where

$$T_B^1 = T_b + \frac{Q_{\text{met}}^1}{\omega_b^1 c_b \rho_b} \text{ and } T_B^2 = T_b + \frac{Q_{\text{met}}^2}{\omega_b^2 c_b \rho_b}.$$

Eq. (8) was solved in **Appendix 2** at the thermal equilibrium. The general solutions are obtained:

$$T_1(r) = c_1 \frac{\cosh(\beta_1 r)}{r} + c_2 \frac{\sinh(\beta_1 r)}{r} + v_1 \frac{\cosh(\beta_1 r)}{r} + v_2 \frac{\sinh(\beta_1 r)}{r} \quad (\text{A2.3})$$

and

$$T_2(r) = c_3 \frac{\cosh(\beta_2 r)}{r} + c_4 \frac{\sinh(\beta_2 r)}{r} + v_3 \frac{\cosh(\beta_2 r)}{r} + v_4 \frac{\sinh(\beta_2 r)}{r} \quad (\text{A2.4})$$

The expressions v_1, v_2, v_3 and v_4 are:

$$v_1 = \frac{1}{\beta_1} \int r [a_1 + b_1 \Phi_1(r)] \sinh(\beta_1 r) dr \text{ and } v_2 = -\frac{1}{\beta_1} \int r [a_1 + b_1 \Phi_1(r)] \cosh(\beta_1 r) dr$$

$$v_3 = \frac{1}{\beta_2} \int r [a_2 + b_2 \Phi_2(r)] \sinh(\beta_2 r) dr \text{ and } v_4 = -\frac{1}{\beta_2} \int r [a_2 + b_2 \Phi_2(r)] \cosh(\beta_2 r) dr$$

with the following notations:

$$a_1 = \beta_1^2 T_B^1 \text{ and } a_2 = \beta_2^2 T_B^2$$

$$b_1 = \frac{\beta_1^2 \bar{P}}{\omega_b^1 c_b \rho_b} \text{ and } b_2 = \frac{\beta_2^2 \bar{P}}{\omega_b^2 c_b \rho_b}$$

$$\beta_1^2 = \frac{\omega_b^1 c_b \rho_b}{k_1}, \beta_2^2 = \frac{\omega_b^2 c_b \rho_b}{k_2}$$

The integration constants: c_1, c_2, c_3, c_4 were computed from the following boundary conditions.

Boundary conditions:

(i) The temperature T_1 is finite at the center ($r \rightarrow 0$) of the geometric structure (**Figure 1**). As a result, the constant c_1 is zero. Therefore, the temperatures: $T_1(r)$ and $T_2(r)$ are (**Appendix 2**):

$$\begin{aligned}
 T_1(r) &= c_2 \frac{\sinh(\beta_1 r)}{r} + v_1 \frac{\cosh(\beta_1 r)}{r} + v_2 \frac{\sinh(\beta_1 r)}{r} \\
 T_2(r) &= c_3 \frac{\cosh(\beta_2 r)}{r} + c_4 \frac{\sinh(\beta_2 r)}{r} + v_3 \frac{\cosh(\beta_2 r)}{r} + v_4 \frac{\sinh(\beta_2 r)}{r}
 \end{aligned}
 \tag{12}$$

(ii) Dirichlet boundary condition was considered on the external surface of the healthy tissue:

$$T_2[r = R_2] = 37^\circ \text{C} \tag{13}$$

(iii) The Newman boundary conditions are considered at malignant—healthy tissue interface. The heat flux coming from the malignant tissue is completely received by the healthy region. The continuity condition of the heat fluxes is imposed at tumor—healthy region interface:

$$k_1 \left[\frac{\partial T_1}{\partial r} \right]_{R_1} = k_2 \left[\frac{\partial T_2}{\partial r} \right]_{R_1} \tag{14}$$

$$T_1[r = R_1] = T_2[r = R_1] \tag{15}$$

3. Results and discussions

In this analysis, the magnetite system with sizes in the range (5–30) nm was considered. The temperature field within a liver tissue was computed for different values of the magnetic field parameters. The malignant and healthy tissues are two concentric domains having the diameter of 20 mm and 100 mm, respectively (**Figure 1**). The temperature values depend strongly on the particle size and magnetic field parameters (frequency and amplitude). The values of the magnetic field parameters: (H_0 and f) verify the criterion of exposure safe and tolerable for the ablation of the whole tumor according with Hergt condition: $H_0 f < 5 \cdot 10^9 \text{Am}^{-1} \text{s}^{-1}$ [7]. Eqs. (1), (3), (5) and (8) were solved in a numerical model using the finite element method (FEM). Their numerical solutions were compared by the previous analytical solutions in the same mathematical conditions.

Figure 2(a) shows the radial dependence of the particle velocity for values of the ferrofluid infusion rate Q_v in the range of 5–30 $\mu\text{l}/\text{min}$. The velocity of the particles on radial direction, decreases with the distance from IS according with the relation (2). The parameter Q_v influences significantly the particles velocities within tissues. **Figure 2(b)** shows the radial dependence of the pressure for the same values of the parameter Q_v . In agreement with relations (4.1) and (4.2), the pressure decreases with the distance from IS. The pressure developed within geometry as result of ferrofluid infusion depends strongly on the Q_v . Higher values of Q_v determines higher values of pressure and faster movements of the particles within tissues. As result of the convection process, the particles move on larger distances or remain in the small vicinity of the IS (at small radial distances). The parameter Q_v and implicitly the pressure generated within tissue by the ferrofluid infusion influences strongly the convection—diffusion—deposition processes of the particles and implicitly the spatial distribution of the particles.

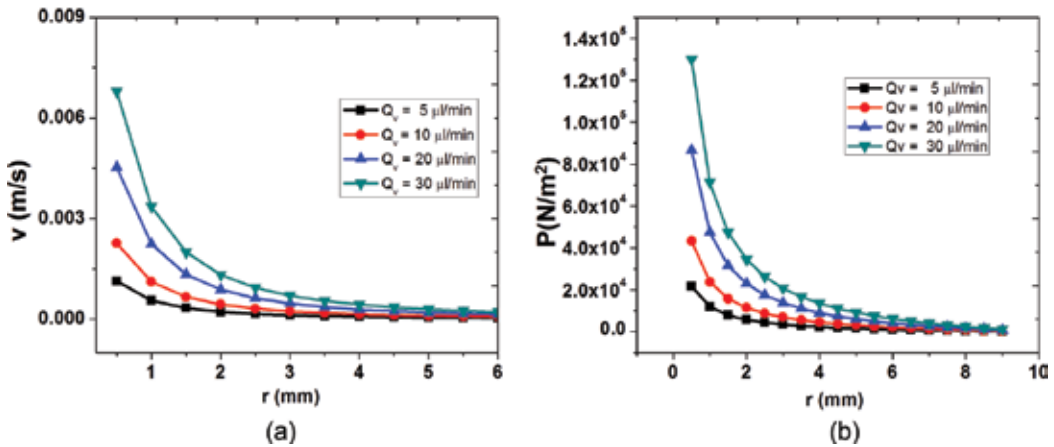


Figure 2. (a) The radial velocity of the magnetite particles within tissues; (b) the pressure on radial direction within malignant tissue.

The radial distribution of the particles and temperature field were analyzed for different values of the parameter Q_v in the range 10–40 $\mu\text{l}/\text{min}$. Figure 3(a) shows the evolution with distance from IS of the volume fraction of the particles $\Phi_1(r)$ within tumor for different values of Q_v . The value of the concentration at IS was $C_{\text{max}} = 10 \text{ mg}/\text{cm}^3$. Consequently, the maximum value of the volume fraction at IS was $\Phi_{\text{max}} = 1.93 \cdot 10^{-3}$. Variation of the parameter Q_v determines different spatial distributions of the particles within the malignant tissue which influences strongly the temperature field (Figure 3(b)). Practically, the temperature values within the malignant tissue can be controlled in the therapeutic range (42–46) $^{\circ}\text{C}$ by using an optimum value of Q_v during the ferrofluid infusion process.

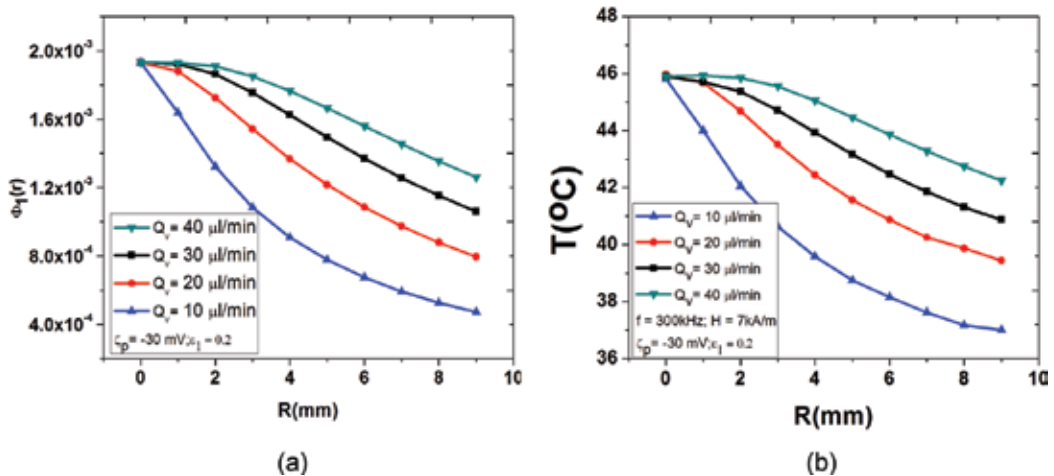


Figure 3. The influence of the parameter Q_v on the radial dependent volume fraction of the particles (a) and the temperature field (b).

The evolution with the parameter Q_v of the deposition rate coefficient of the particles k_f^1 was studied for different porosities of the malignant tissue, in order to understand the influence of this parameter on the deposition process of the particles on the solid porous matrix (**Figure 4**). The coefficient k_f^1 decreases with the increase of the value of the parameter Q_v . Also the deposition of the particles decreases with the increase of the tissue porosity. The particles with high velocities (high values of Q_v) as a result of the high pressure gradient have no capability to remain deposited on the solid matrix.

The repulsive electrostatic double layer forces influences the particle deposition process and implicitly the spatial distribution of the particles. Temperature field depends strongly on the particle zeta potential ζ_p . **Figure 5(a)** shows the evolution with radial distance from IS of the volume fraction of the particles $\Phi_1(r)$ for different values of the particle zeta potential $\zeta_p = -10$ to -40 mV. As a result of the strong repulsive electrostatic double layer forces, a number of the particles are deposited in the solid structure of tissue. This effect influences significantly the spatial distribution of the temperature (on radial direction) as **Figure 5(b)** shows. As a consequence, the temperature gradients become smaller in the case of smaller repulsive electrostatic double layer forces. The repulsive electrostatic interactions (as a result of the repulsive electrostatic double layer (EDL) forces) influence strongly the mass concentration of the particles and the spatial temperature field. The particle zeta potential ζ_p can be controlled in the ferrofluid (as liquid medium) due to the ionic conditions measured by pH and ionic strength.

The values of the magnetic field parameters: (H_0 and f) are essential in the optimization of the Magnetic Hyperthermia therapy. In the following, the increase of the temperature on the radial direction with the frequency of the magnetic field was followed.

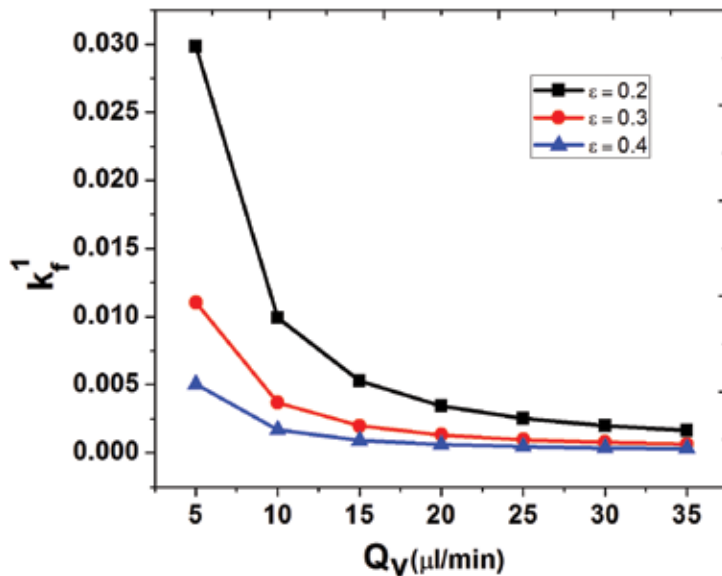


Figure 4. The evolution with the parameter Q_v of deposition rate coefficient of the particles, k_f^1 within the malignant tissue.

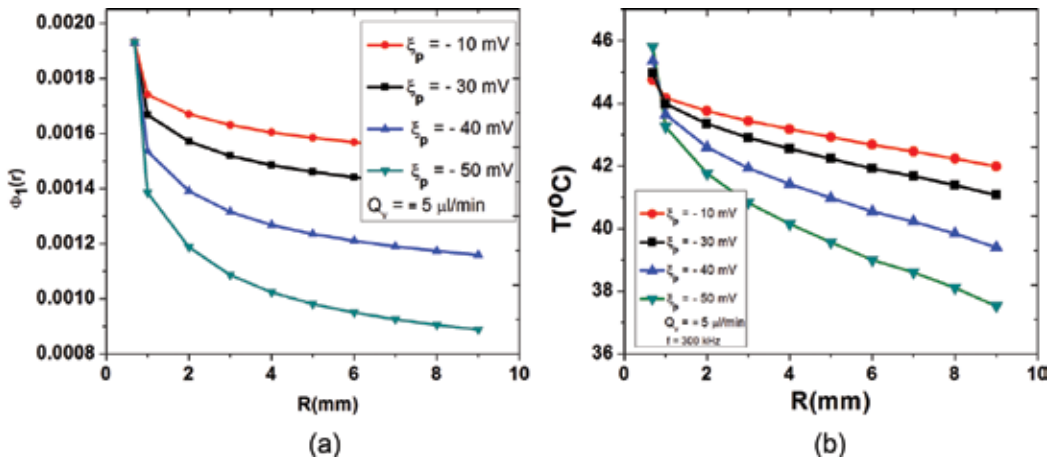


Figure 5. The influence of the parameter—zeta potential ξ_p on the radial dependent volume fraction of the particles (a) and the temperature field (b).

Figure 6(a) and (b) show a 3D and 2D view of the radial temperature field within the malignant tissue for different values of the frequency of the magnetic field. It was considered a small value of both Q_v and the tissue porosity in order to analyze a temperature field with strong non-uniformity (and implicitly high thermal gradients).

Figure 7(a) shows the values of the main parameters Q_v and f which determines the same temperature on the radial direction. Figure 7(b) shows the isothermal surfaces for different values of values of the main parameters Q_v and f .

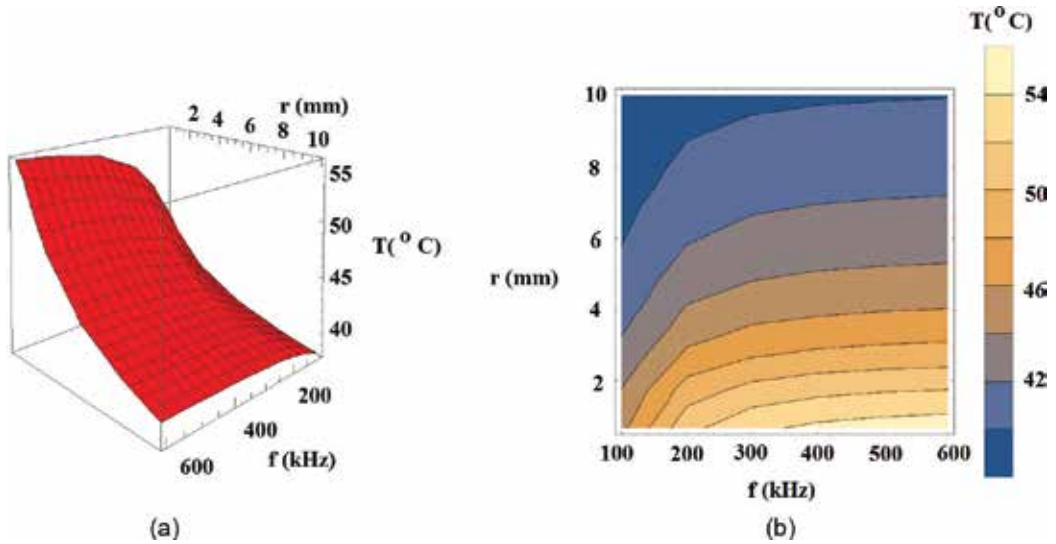


Figure 6. The evolution with the frequency of the magnetic field of the temperature field on radial direction $Q_v = 10 \mu\text{l}/\text{min}$; $\epsilon_1 = 0.2$ and $\xi_p = -30 \text{ mV}$. (a) 3D view and (b) 2D view.

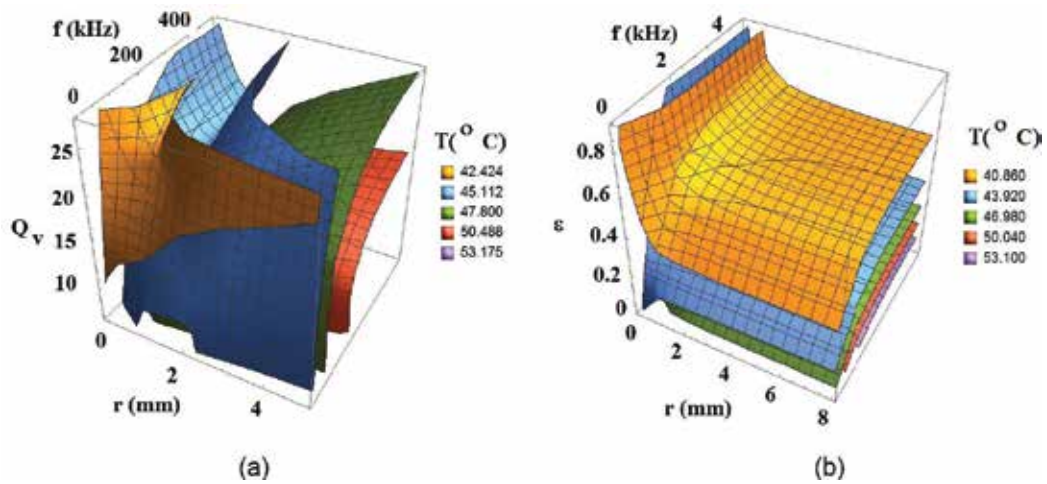


Figure 7. The isothermal surfaces for different parameters. (a) Isothermal surfaces for different values of the parameter Q_v and frequency of the magnetic field f and (b) isothermal surfaces for different values of the tissue porosity and frequency of the magnetic field f .

The analytical temperature was compared with the numerical results given by the finite element method (FEM) in Comsol Multiphysics using the same parameters in the same mathematical conditions. Good agreements were found between the predictions of the analytical model and numerical results.

The simulations allow (i) the optimization of the main parameters which influences strongly the heating of the tumor in the therapeutic temperature range and (ii) provide better temperature control through treatment preplanning.

4. Conclusion

The model developed in this paper analyzes the essential role of the ferrofluid infusion rate in the radial MNP distribution after their injection within a malignant tissue. Analytical correlations between the following parameters: (i) the particle velocity, (ii) the pressure developed in geometry during the ferrofluid infusion and (iii) the particle concentration were done in order to understand and predicts the temperature field within tissues when an external magnetic field is applied. The temperature field is concentrated within the malignant tissue. The temperature on the tumor border (approximately 38–39 $^{\circ}\text{C}$) not affects the healthy tissue.

The thermal energy deposited within the malignant tissue provides from the MNP distributed as result of convection-diffusion-deposition of the particles after their injection inside tissue. The ferrofluid infusion rate influences significantly the radial distribution of the particles and consequently the temperature field.

The temperature field within the malignant tissues can be controlled by the control of the ferrofluid infusion rate Q_v during the infusion process. The particles having higher velocity moves

on larger distances on radial direction from the injection site within tumor. As a result, the particles which not remain in the vicinity of the injection site are distributed in the tumor volume. This important effect determines a temperature field with small temperature gradients. The model developed in this paper can be used as a planning tool to compute the temperature field for different parameters.

Appendix 1

2.1. Convection-Diffusion-Deposition of the particles

1) The computation of the deposition rate coefficients of the particles, k_f^i

The mean deposition rate coefficients of the particles are expressed by the following relation [10–16]:

$$k_f^i = \frac{3(1 - \varepsilon_i)}{2 \varepsilon_i d_c} \eta_s^i v_r, \quad (i = 1, 2) \tag{A1.1}$$

The porous media contains the spherical collector cells with diameters d_c ranged from 0.05 to 0.50 mm. The coefficients k_f^i depend on the particle diameter D , collector diameter d_c , porosities of the malignant and healthy tissues ε_p and the radial velocity v_r . The *collector efficiency* η_s^i describes the ratio of the particles captured by the solid surface to those brought into a unit structural cell of the porous medium [8, 13]. This coefficient is given by the expression:

$$\eta_s^i = \alpha^i \eta_0^i \tag{A1.2}$$

The *single collector contact efficiency* η_0^i is the fraction of the particles brought to the collector surface by the Brownian diffusion, interception and/or gravitational sedimentation. This coefficient was computed by N. Tufenkji and M. Elimelech considering the superposition of the effects developed by the hydrodynamic forces, van der Waals interactions and gravity effect [13]:

$$\eta_0^i = \underbrace{\eta_D^i}_{\substack{\text{transport by} \\ \text{diffusion}}} + \underbrace{\eta_I^i}_{\substack{\text{transport by} \\ \text{interception}}} + \underbrace{\eta_G}_{\substack{\text{transport by} \\ \text{gravitation}}}$$

with:

$$\eta_D^i = 2.4 (A_s^i)^{1/3} N_R^{-0.081} (N_{Pe})^{-0.715} N_{vdW}^{0.052}; \quad \eta_I^i = 0.55 A_s^i N_R^{1.675} (N_A)^{0.125};$$

$$\eta_G = 0.22 N_R^{-0.24} (N_G)^{1.11} N_{vdW}^{0.053}$$

The following non dimensional coefficients are defined [13]:

1. $N_R = \frac{D}{d_c}$ —the aspect ratio;
2. $N_{Pe} = \frac{U d_c}{D_w}$ —Peclet number which defines the ratio of the convective transport to diffusive transport;
3. $N_{vdw} = \frac{A}{kT}$ —van der Waals number—the ratio of van der Waals interaction energy to the particle's energy;
4. $N_A = \frac{A}{12 \pi \mu R^2 U}$ —attraction number—combined influence of van der Waals attraction forces and fluid velocity on particle deposition rate due to interception;
5. $N_G = \frac{2 R^2 (\rho_{MNP} - \rho_{ferro}) g}{9 \mu U}$ —gravity number—the ratio of Stokes particle settling velocity to approach velocity of the fluid;
6. $A_s^i = \frac{2(1-\gamma_i^5)}{2-3\gamma_i+3\gamma_i^5-2\gamma_i^6}$; $\gamma_i = (1-\varepsilon_i)^{1/3}$ —porosity dependent parameter of Happel's model.

Deposition of the particles on the pore wall is influenced by the electrostatic repulsive forces. The *attachment (collision) efficiency coefficient* (filter coefficient) α^i ($i = 1, 2$) represents the fractional reduction in deposition rate of the particles due to the presence of the electrostatic repulsive energies [15]. Bai and Tien (1999), Chang and Chan [16] computed the expression of α^i . The analytical expression was compared successfully with the experimental data. They found the following correlation equation:

$$\alpha^i = \exp\left[\frac{1}{2}(\ln(\alpha_{C-C}^i) + \ln(\alpha_{B-T}^i))\right], (i = 1, 2). \tag{A1.3}$$

The expression of α_{C-C}^i is given by Chang and Chan [15, 16]:

$$\alpha_{C-C}^i = 0.024 (N_{DL})^{0.969} (N_{E1})^{-0.423} (N_{E2})^{2.880} (N_{LO})^{1.5} + 3.176 (A_s^i)^{1/3} (N_R)^{-0.081} (N_{Pe}^i)^{-0.715} (N_{LO})^{2.687} + 0.222 A_s^i (N_R)^{3.041} (N_{Pe})^{-0.514} (N_{LO})^{0.125} + (N_R)^{-0.24} (N_G)^{1.11} (N_{LO}) \tag{A1.4}$$

The expression of α_{B-T}^i is given by Bai and Tien [15, 16]:

$$\alpha_{B-T}^i = 2.527 \cdot 10^{-3} (N_{LO})^{0.7031} (N_{E1})^{-0.3121} (N_{E2})^{3.5111} (N_{DL})^{1.352} \tag{A1.5}$$

In the absence of the electrostatic double layer forces, α^i become 1. The nondimensional coefficients: N_{LO} , N_{E1} , N_{E2} and N_{DL} have the following expressions:

7. $N_{LO} = \frac{4A}{36 \pi \mu R^2 U}$ is London number;
8. $N_{E1} = \frac{\varepsilon_r \varepsilon_0 (\zeta_p^2 + \zeta_c^2)}{6 \pi \mu R U}$ is the first electrokinetic parameter;

9. $N_{E2} = \frac{2 \zeta_p \zeta_c}{\zeta_p^2 + \zeta_c^2}$ is the second electrokinetic parameter;

10. $N_{DL} = 2 \kappa R$ is the double layer force parameter;

κ is Debye length for the colloidal suspension; ζ_p is particle zeta potential; ζ_c is collector zeta potential (**Table 1**) and $u = \frac{Q_v}{S_{needle}}$ is the ferrofluid velocity at the top of the needle.

The repulsive electrostatic double layer (EDL) forces appear in the liquid medium due to the ionic conditions measured by pH and ionic strength.

2. The computation of the MNP concentrations $C_i = C_i(r)$ as solution of Eq. (5)

At equilibrium, in the steady-state: $\frac{\partial C_i}{\partial t} = 0$ and Eq. (5) becomes:

$$\nabla \cdot (\vec{v} C_i) = \nabla \cdot (D_i^* \nabla C_i) - k_f^i C_i, \quad (\text{A1.6})$$

where the deposition rate coefficients of the particles k_f^i are given by the relations (A1.1):

$$\frac{1}{r^2} \frac{\partial}{\partial r} (v_r r^2 C_i) - \frac{1}{r^2} \frac{\partial}{\partial r} \cdot \left(D_i^* r^2 \frac{\partial C_i}{\partial r} \right) = -k_f^i C_i \quad (\text{A1.7})$$

$$\frac{\partial}{\partial r} \left(v_r r^2 C_i - D_i^* r^2 \frac{\partial C_i}{\partial r} \right) = M_i C_i \quad (\text{A1.8})$$

$$\frac{\partial}{\partial r} \left\{ C_i r^2 \left(\frac{A_i}{r^2} - \frac{1}{C_i} \frac{\partial C_i}{\partial r} \right) \right\} = m_i C_i \quad (\text{A1.9})$$

with the following constants:

$$M_i = -\frac{3(1 - \varepsilon_i)}{2 \varepsilon_i d_c} \eta_s^i B; m_i = \frac{M_i}{D_i^*}; A_i = \frac{B}{D_i^*}; B = \frac{Q_v}{\pi} \quad (\text{A1.10})$$

Solution of Eq. (A1.9) has the following form:

$$C_i(r) = C_0^i(r) \exp \left[-\frac{A_i}{r} \right] \quad (\text{A1.11})$$

Considering $C_0^i = C_0^i(r)$, Eq. (A1.10) can be written as:

$$-\frac{\partial}{\partial r} \left\{ \left(r^2 \exp \left[-\frac{A_i}{r} \right] \right) \frac{\partial C_0^i}{\partial r} \right\} = m_i C_0^i \exp \left[-\frac{A_i}{r} \right] \quad (\text{A1.12})$$

or:

$$r^2 \left(\frac{\partial^2 C_0^i}{\partial r^2} \right) + (2r + A_i) \left(\frac{\partial C_0^i}{\partial r} \right) + m_i C_0^i = 0$$

which is equivalent with

$$\frac{\partial}{\partial r} \left[r^2 \frac{\partial C_0^i}{\partial r} \right] + A_i \left(\frac{\partial C_0^i}{\partial r} \right) + m_i C_0^i = 0 \tag{A1.13}$$

The solutions of Eq. (A1.13) are given by the following expressions:

$$C_0^i(r) = \frac{(\text{const1}) e^{\frac{A_i}{2r}} \text{Bessel I} \left[\sqrt{1-4m_i}, \frac{A_i}{2r} \right]}{\sqrt{r}} + \frac{(\text{const2}) e^{\frac{A_i}{2r}} \text{Bessel K} \left[\sqrt{1-4m_i}, \frac{A_i}{2r} \right]}{\sqrt{r}} \tag{A1.14}$$

$\text{Bessel I} \left[\sqrt{1-4m_i}, \frac{A_i}{2r} \right]$ and $\text{Bessel K} \left[\sqrt{1-4m_i}, \frac{A_i}{2r} \right]$ are modified Bessel functions I and K of the order $\sqrt{1-4m_i}$. The expressions $\frac{A_i}{2r}$ are the variables of these functions. The general solutions of Eq. (5) are computed using the expressions (A1.14) in the expression (A1.11):

$$C_i(r) = \left(\frac{e^{\frac{A_i}{2r}}}{\sqrt{r}} \right) \left((\text{const1})_i \text{Bessel I} \left[\sqrt{1-4m_i}, \frac{A_i}{2r} \right] + (\text{const2})_i \text{Bessel K} \left[\sqrt{1-4m_i}, \frac{A_i}{2r} \right] \right) \tag{A1.15}$$

$(\text{const1})_i$ and $(\text{const2})_i$ are the four integration constants which are determined from the following four boundary conditions:

- i. $C_2 = 0$ on the external boundary of the geometry ($r = R_2$);
- ii. Neumann boundary condition at the all inner interfaces;

$$C_1(r = R_1) = C_2(r = R_1)$$

$$D_1^* \frac{\partial C_1}{\partial r} \Big|_{r=R_1} = D_2^* \frac{\partial C_2}{\partial r} \Big|_{r=R_1}$$

- iii. at the injection site (IS), at the top of the needle ($r = r_0$) the concentration has the particular expression $C_1 = C_{\max}$.

The constants $(\text{const1})_i$ and $(\text{const2})_i$ were computed in the Wolfram Mathematica 10 software.

Appendix 2

2.2. The temperature model

At the thermal equilibrium, Eq. (7) are:

$$\frac{1}{r^2} \frac{\partial}{\partial r} \left[r^2 \frac{\partial T_i}{\partial r} \right] + \beta_i^2 (T_e^i(r) - T_i) = 0 \text{ or } \frac{\partial^2 T_i}{\partial r^2} + \frac{2}{r} \frac{\partial T_i}{\partial r} - \beta_i^2 T_i + \beta_i^2 T_e^i(r) = 0 \tag{A2.1}$$

Using substitution $R_i = rT_i$, these equations become:

$$\frac{d^2 R_i}{dr^2} - \beta_i^2 R_i + \beta_i^2 r T_e^i(r) = 0 \quad (\text{A2.2})$$

with solutions

$$R_1(r) = c_1 \cosh(\beta_1 r) + c_2 \sinh(\beta_1 r) + v_1 \cosh(\beta_1 r) + v_2 \sinh(\beta_1 r)$$

$$R_2(r) = c_3 \cosh(\beta_2 r) + c_4 \sinh(\beta_2 r) + v_3 \cosh(\beta_2 r) + v_4 \sinh(\beta_2 r)$$

where

$$v_1 = \frac{1}{\beta_1} \int r [a_1 + b_1 \Phi_1(r)] \sinh(\beta_1 r) dr \text{ and } v_2 = -\frac{1}{\beta_1} \int r [a_1 + b_1 \Phi_1(r)] \cosh(\beta_1 r) dr$$

$$v_3 = \frac{1}{\beta_2} \int r [a_2 + b_2 \Phi_2(r)] \sinh(\beta_2 r) dr \text{ and } v_4 = -\frac{1}{\beta_2} \int r [a_2 + b_2 \Phi_2(r)] \cosh(\beta_2 r) dr$$

These expressions contain the following notations:

$$a_1 = \beta_1^2 T_B^1 \text{ and } a_2 = \beta_2^2 T_B^2$$

$$b_1 = \frac{\beta_1^2 \bar{P}}{\omega_b^1 c_b \rho_b} \text{ and } b_2 = \frac{\beta_2^2 \bar{P}}{\omega_b^2 c_b \rho_b}$$

The solutions of Eq. (A2.1) are:

$$T_1(r) = c_1 \frac{\cosh(\beta_1 r)}{r} + c_2 \frac{\sinh(\beta_1 r)}{r} + v_1 \frac{\cosh(\beta_1 r)}{r} + v_2 \frac{\sinh(\beta_1 r)}{r} \quad (\text{A2.3})$$

and

$$T_2(r) = c_3 \frac{\cosh(\beta_2 r)}{r} + c_4 \frac{\sinh(\beta_2 r)}{r} + v_3 \frac{\cosh(\beta_2 r)}{r} + v_4 \frac{\sinh(\beta_2 r)}{r} \quad (\text{A2.4})$$

c_1, c_2, c_3, c_4 are the integration constants and β_i ($i = 1, 2$) have the expressions:

$$\beta_1^2 = \frac{\omega_b^1 c_b \rho_b}{k_1}; \beta_2^2 = \frac{\omega_b^2 c_b \rho_b}{k_2}$$

Author details

Iordana Astefanoaei* and Alexandru Stancu

*Address all correspondence to: iordana@uaic.ro

Faculty of Physics, Alexandru Ioan Cuza University of Iasi, Romania

References

- [1] Obaidat IM, Issa B, Haik Y. *Nanomaterials*. 2015;**5**:63
- [2] Lahonian M, Golneshan AA. *IEEE Transactions on Nanobioscience*. 2011;**10**:4
- [3] Carrey J, Mehdaoui B, Respaud M. *Journal of Applied Physics*. 2011;**109**:083921
- [4] Salloum M, Ma RH, Weaks D, Zhu L. *International Journal of Hyperthermia*. 2008;**24**:4
- [5] Salloum M, Ma R, Zhu L. *International Journal of Hyperthermia*. 2009;**25**:4
- [6] Astefanoaei I, Stancu A, Chiriac H. *AIP Conference Proceedings*. 2017;**1796**:040006
- [7] Rosensweig RE. *Journal of Magnetism and Magnetic Materials*. 2002;**252**:370-374
- [8] Su D, Ma R, Salloum M, Zhu L. *Medical & Biological Engineering & Computing*. 2010;**48**:7
- [9] Vafai K, editor. *Handbook of Porous Media*. 2nd ed. Singapore: CRC Taylor & Francis
- [10] Morrison PF, Laske DW, Bobo H, Oldfield EH, Dedrick RL. *American Journal of Physiology. Regulatory, Integrative and Comparative Physiology*. 1994;**266**:21
- [11] Satterfield CN. *Mass Transport in Heterogeneous Catalysis*. Cambridge MA: MIT Press; 1970. 5 p
- [12] Ramanujan S, Pluen A, McKee TD, Brown EB, Boucher Y, Jain RK. *Biophysical Journal*. 2002;**83**:3-10
- [13] Tufenkji N, Elimelech M. *Environmental Science & Technology*. 2004;**38**:20
- [14] Rajagopalan R, Tien C. *AIChE Journal*. 1976;**22**:3
- [15] Tien C. *Principles of Filtration*. Amsterdam: Elsevier; 2012
- [16] Chang Y-I, Chan H-C. *AIChE Journal*. 2008;**54**:5
- [17] Golneshan AA, Lahonian M. *International Journal of Hyperthermia*. 2011;**27**:3

Vibration Modeling and Numerical Acoustics

Direct and Hybrid Aeroacoustic Simulations Around a Rectangular Cylinder

Hiroshi Yokoyama and Akiyoshi Iida

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.70810>

Abstract

Aeroacoustic simulations are divided into hybrid and direct simulations. In this chapter, the effects of freestream Mach number on flow and acoustic fields around a two-dimensional square cylinder in a uniform flow are focused on using direct and hybrid simulations of flow and acoustic fields are performed. These results indicate the effectiveness and limit of the hybrid simulations. The Mach number M is varied from 0.2 to 0.6. The propagation angle of the acoustic waves for a high Mach number such as $M = 0.6$ greatly differs from that predicted by modified Curle's equation, which assumes the scattered sound to be dominant and takes the Doppler effects into consideration. This is because the acoustic field is affected by the direct sound, which is generated by quadrupoles in the original Curle's equation. To clarify the effects of the direct sound on the acoustic field, the scattered and direct sounds are decomposed. The results show that the direct sound is too intense to neglect for $M \geq 0.4$. Moreover, acoustic simulations are performed using the Lighthill's acoustic sources.

Keywords: aeroacoustics, aeolian tone, direct simulation, acoustic analogy, Lighthill's equation

1. Introduction

The sound generated by a cylinder in a uniform flow is known as the aeolian tone. This sound is often radiated from flows around a cylinder. Strouhal [1] found that the frequency of the tone is identical to the vortex shedding frequency. Lighthill [2] derived the nonhomogeneous wave equation as shown in Eqs. (1) and (2) from the compressive Navier-Stokes equations.

$$\frac{\partial^2}{\partial t^2} \rho - a_0^2 \frac{\partial^2}{\partial x_i \partial x_j} \rho = \frac{\partial^2 T_{ij}}{\partial x_i \partial x_j}, \quad (1)$$

where ρ is the density, a_0 is the freestream sound speed, and the tensor T_{ij} is defined by:

$$T_{ij} = \rho v_i v_j + \delta_{ij}((p - p_0) - a_0^2(\rho - \rho_0)) - \tau_{ij}, \quad (2)$$

where v_i is velocity, p is pressure, the tensor τ is the viscous stress tensor, and $\partial^2 T_{ij} / \partial x_i \partial x_j$ is referred to as Lighthill's acoustic source. The first, second, and third terms of the Lighthill's acoustic source are related to the momentum, entropy, and viscosity, respectively. For aerodynamic sound around a body in a fluid stream of a low Mach number, Curle [3] has shown through analytical solution of Lighthill's equation [2] that the surface pressure fluctuations around the body lead to a dipole sound field. Investigations such as those by Gerrard [4] and Phillips [5] have experimentally confirmed that the acoustic field around a circular cylinder has directivity normal to the fluid stream and is closely related to the fluctuations of the lift force.

Recently, many investigations using numerical simulations have been performed, for instance, Inoue and Hatakeyama [6], Gloerfelt et al. [7], and Liow et al. [8]. Inoue and Hatakeyama [6] modified the Curle's solution considering the Doppler effects and showed that the acoustic fields predicted by the proposed equation agree well with those predicted by their direct simulations for a low freestream Mach number $M \leq 0.3$. Gloerfelt et al. [7] performed incompressible flow computations and acoustic computations on the basis of Lighthill's acoustic analogy [2] for the flow around a circular cylinder. The role of the acoustic scattering on the cylinder in the mechanism of the sound generation was investigated for $M = 0.12$ and the Reynolds number based on the diameter $Re_d \approx 1.1 \times 10^5$. Here, the entropy (second) and viscous (third) terms in Lighthill's acoustic source Eq. (2) were neglected. This is reasonable for the flow of such a high Reynolds number and a low Mach number [9]. Liow et al. [8] also performed the incompressible flow simulations and acoustic simulations for a flow around an elongated rectangular cylinder with $M \leq 0.2$. The acoustic computations are based on the Powell's theory [10], where the acoustic sources approximately correspond to the momentum (the first term) of Lighthill's acoustic source. The effects of the drag force on the acoustic field were clarified.

Despite many investigations into the aeolian sound around a cylinder, little attention has been given to flows around a cylinder with a high Mach number $M > 0.3$. For such a high Mach number, the effects of the Mach number on the flow and acoustic fields around a cylinder have not been clarified. Also, it is currently unknown whether the contribution of the second and third terms in Lighthill's acoustic source to the acoustic field can be neglected for such a high Mach number. In high-speed jets such as $M = 0.9$ – 2.0 , it has been clarified that the second term needs to be taken into consideration [11].

In the present chapter, aerodynamic sound radiated from a two-dimensional square cylinder in a freestream is investigated. The flow field around a square cylinder has been investigated by many researchers [12–14]. However, little is known about the acoustic field. The hybrid and direct simulations of flow and acoustic fields are introduced. The freestream Mach number on the flow and acoustic fields are focused on. The Mach number is varied from 0.2 to 0.6. Moreover, the contributions of each term of Lighthill's acoustic source to the acoustic field are focused on. To do this, the acoustic simulations are also performed using the Lighthill's acoustic sources computed by the direct simulations. This method for predicting the acoustic field using the acoustic simulation is referred to as the hybrid simulation in this chapter.

2. Numerical methods

2.1. Flow configurations

The flow around a two-dimensional square cylinder, as shown in **Figure 1**, is investigated. To clarify the effects of the freestream Mach number on flow and acoustic fields, the computations are performed for $M = 0.2, 0.3, 0.4, 0.5,$ and 0.6 . The Reynolds number based on the freestream velocity and the side length of the cylinder is set to 150, where the three-dimensional instability does not occur [13]. Here, the two-dimensional phenomena related to the vortex shedding from the cylinder are focused on.

The fluid was assumed to be standard air, where Sutherland's formula can be applied for the viscosity coefficient. The specific heat C was assumed to be $1004 \text{ J kg}^{-1} \text{ K}^{-1}$ and that the Prandtl number Pr was 0.72.

2.2. Direct simulation

2.2.1. Governing equations and finite difference formulation

Both flow and acoustic fields are solved by the two-dimensional compressible Navier-Stokes equations in a conservative form, which is written as:

$$Q_t + (E - E_v)_x + (F - F_v)_y = 0, \quad (3)$$

where Q is the vector of the conservative variables, E and F are the inviscid fluxes, and E_v and F_v are the viscous fluxes. The spatial derivatives and time integration were evaluated by the sixth-order accurate compact finite difference scheme [15] and a third-order accurate Runge-Kutta method. To suppress the numerical instabilities associated with the central differencing in the compact scheme, we use a tenth-order accurate spatial filter shown below:

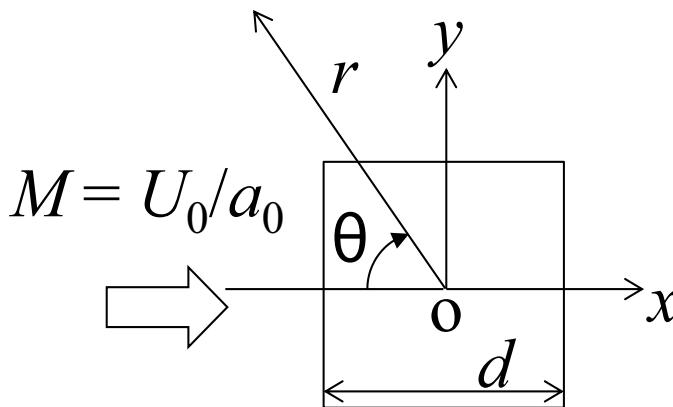


Figure 1. Configurations for flow around a two-dimensional square cylinder.

$$\alpha_f \hat{\varphi}_{i-1} + \hat{\varphi}_i + \alpha_f \hat{\varphi}_{i+1} = \sum_{n=0}^5 \frac{a_n}{2} (\varphi_{i+n} + \varphi_{i-n}), \quad (4)$$

where φ is a conservative quantity and $\hat{\varphi}$ is the filtered quantity. The coefficients a_n are the same as the values used by Gaitonde and Visbal [16], and the parameter α_f is 0.47.

2.2.2. Computational grids and boundary conditions

Figure 2 shows the computational domain and boundary conditions. The coordinates originate from the center of the cylinder. Generally, the nonreflecting boundary conditions based on the characteristic wave relations [17–19] are used at the inflow, upper, and outflow boundaries along with a buffer region. The role of the buffer region is similar to that of the “sponge region” of Colonius et al. [20]. At the wall, the nonslip and adiabatic boundary conditions are used.

For all the cases of $M = 0.2$ – 0.6 , the same grids are used. The computational domain is divided into three regions of different grid spacings as shown in **Figure 2**: a vortex region $[-4.0 \leq x/d \leq 30.0, -4.0 \leq y/d \leq 4.0]$, a sound region $[-70.0 \leq x/d < -4.0, 30.0 < x/d \leq 70.0, -70.0 \leq y/d < -4.0, 4.0 < y/d \leq 70.0]$, and a buffer region $[-500.0 \leq x/d < -70.0, 70.0 < x/d \leq 500.0, -500.0 \leq y/d < -70.0, 70.0 < y/d \leq 500.0]$.

The spacing in the vortex region is prescribed to be fine enough to analyze the separated shear layer and the vortical structures in the wake of the cylinder. **Figure 3** shows the computational grid near the cylinder. The spacing adjacent to the cylinder surface is $\Delta x_{\min}/d$ and $\Delta y_{\min}/d = 0.0025$. With this grid distribution, the number of grid points within the separated shear layer for $Re = 150$ is 22 in the x, y direction (the thickness of the separated shear layer was estimated by $\delta/d \sim 1/Re^{0.5}$ and 0.08 for $Re = 150$ like the circular cylinder [6]), and the separated shear layer can be sufficiently captured. In the whole vortex region, $\Delta x/d$ and $\Delta y/d$ are less than 0.2, where the

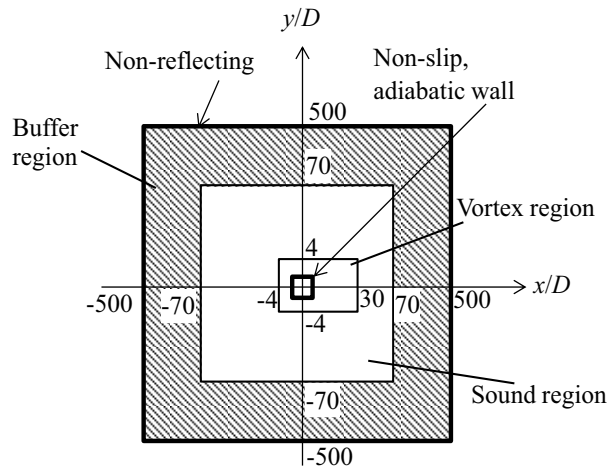


Figure 2. Computational domain and boundary conditions for direct simulations.

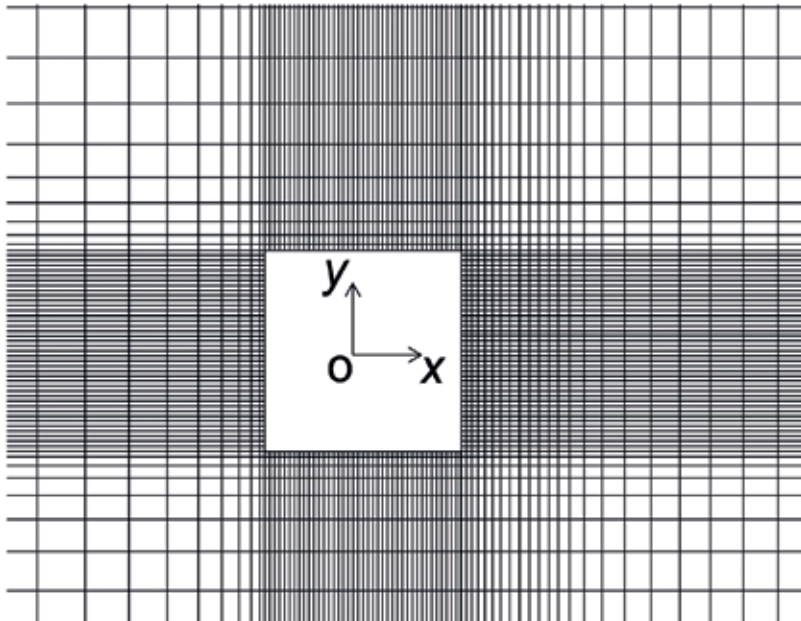


Figure 3. Computational grid near cylinder for direct simulations. Every 10th grid line is shown for clarity.

number of grid points within a shed vortex is about 15 in the x, y direction (the size of the vortex is estimated to be about $3d$ by the spacing of the local maxima of vorticity in the computational results) and the vortices were sufficiently analyzed.

In the sound region, the spacing is prescribed to be larger than that in the vortex region but still fine enough to capture the acoustic waves. The spacings are $\Delta x/d, \Delta y/d \leq 0.23$ except for the downstream region of the cylinder. In the downstream region, the largest spacings are $\Delta x/d$ and $\Delta y/d = 0.46$. This is because the acoustic wavelength becomes longer than that in the upstream region due to the Doppler effects. In the whole sound region, more than 20 grid points are used per one acoustic wavelength of the tonal sound at the frequency of the vortex shedding, and the acoustic waves are sufficiently captured.

After many preliminary tests, grid- and domain-size independence has been established for the solutions presented in this chapter.

2.3. Hybrid simulation

2.3.1. Governing equations and discretization formulation

The two-dimensional Lighthill's equation [Eqs. (1) and (2)] is solved based on the wave equation. Here, the open-source software, FrontFlow/blue-ACOUSTICS, was used. Here, the acoustic simulations are performed in a frequency domain using finite-element methods. A component perturbed at the frequency f of quantify g_f can be written as:

$$g_f = \tilde{g}_f(\mathbf{x})e^{2\pi ft}. \quad (5)$$

Using Eq. (5), Lighthill's equation can be written as:

$$\frac{\partial^2 \tilde{\rho}_f}{\partial x_i \partial x_j} + k^2 \tilde{\rho}_f = -\frac{1}{a_0^2} \frac{\partial^2 \tilde{T}_{ij,f}}{\partial x_i \partial x_j}, \quad (6)$$

where $k = 2\pi f/a_0$ is the wavenumber. The right-hand side of Eq. (6) is computed by the results of the direct simulation in the present chapter. Also, acoustic waves at the frequency of the vortex shedding are focused on. To minimize the spurious errors, the computed acoustic sources are reduced smoothly to zero near the outflow boundary of the acoustic simulation by using the filter. **Figure 4** shows the computational domain for the acoustic simulations and the outer shape of computational domain is circular with the cylinder at its center, and the radius is $100D$. The above-mentioned filter is defined as:

$$\hat{A} = A \times G(r - r_0), \quad (7)$$

$$A = \partial^2 \frac{\tilde{T}_{ij,f}}{\partial x_i \partial x_j}, \quad (8)$$

$$G(r_d) = \begin{cases} \frac{1}{2} \left(1 + \cos \frac{r_d}{L} \pi \right) & (80.0 < r \leq 100.0) \\ 1.0 & (r \leq 80.0) \end{cases}, \quad (9)$$

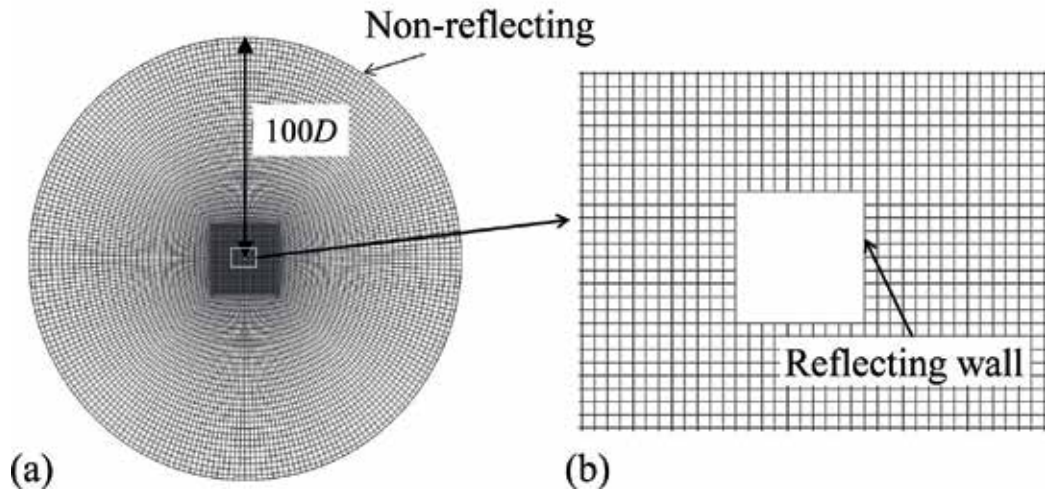


Figure 4. Computational grids and boundary conditions for acoustic simulations. (a) Overall grids (every fifth grid line is shown for clarity) and (b) grids near cylinder (every grid line is shown).

$$r_0/d = 80.0, L/d = 20, \tag{10}$$

where r is the distance from the center of the cylinder.

2.3.2. Computational grids and boundary conditions

Figure 4 shows the computational grid for the acoustic simulations. The spacing adjacent to the cylinder surface is $\Delta x_{\min}/d$ and $\Delta y_{\min}/d = 0.1$. In the whole domain, the grid spacing is less than 0.52, and more than 20 grid points are used per acoustic wavelength. The preliminary computations confirmed that the acoustic waves are sufficiently analyzed with these grid resolutions.

The reflecting conditions are adopted on the cylinder wall. On the other boundaries, the nonreflecting boundary conditions are adopted.

3. Validation of computational methods

3.1. Validation of direct simulations

Figure 5 shows the Strouhal number of vortex shedding predicted by the present direct simulations. The Strouhal number St is the frequency nondimensionalized by the freestream velocity U_0 and the side length of the cylinder d . The present results are compared with the results of the past incompressible simulation ($St = 0.155$) [13] and those of the past experiment ($St = 0.162$) [14] for the same Reynolds number. The flow condition of the experiment is approximately incompressible. The present Strouhal numbers for all the Mach numbers are slightly lower than those in past results. The present computational results show that the Strouhal number becomes lower as the freestream Mach number becomes higher. This is

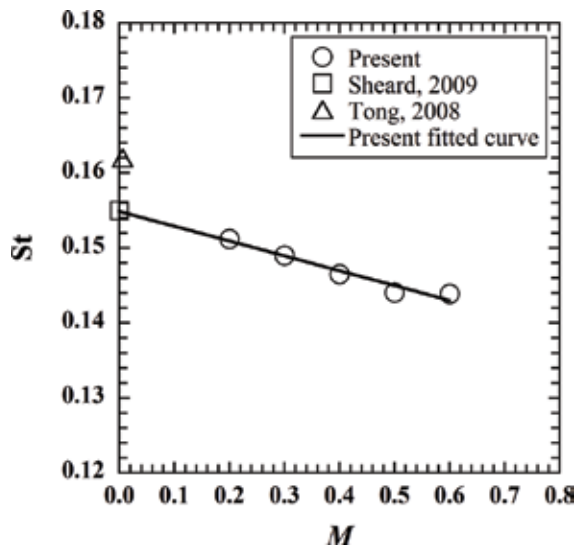


Figure 5. Effects of Mach number on frequency of vortex shedding.

related to the variation of vortices with the variation of the Mach number as discussed in detail in Section 4.1. Also, based on the present results, the extrapolated Strouhal number at $M = 0.0$ is 0.155. This value agrees well with the past computational data [13], although it is not clear why the past experimental value [14] is slightly higher. Consequently, the present direct simulations are confirmed to be validated.

3.2. Validation of hybrid simulation

Figure 6 shows the polar plots of the sound pressure levels at $r/d = 30.0$ predicted by direct and hybrid simulations for $M = 0.4$. The acoustic field by the hybrid simulation is approximately in good agreement with that by the direct simulation. It has been confirmed that the two fields also agree for other Mach numbers such as $M = 0.2$ and 0.6 . The above-mentioned methods of hybrid simulation are clarified to be validated.

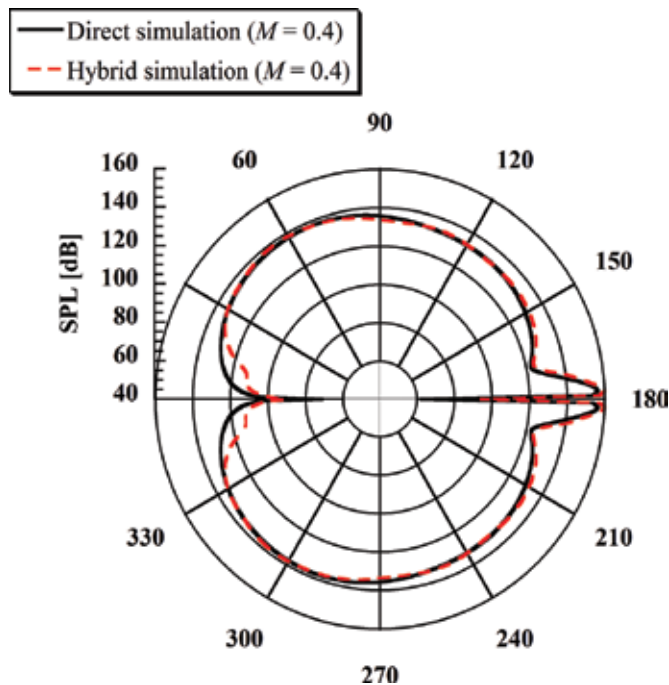


Figure 6. Polar plots of sound pressure levels by direct and hybrid simulations at $r/d = 30.0$ for $M = 0.4$.

4. Results and discussion

4.1. Flow fields

Figure 7 shows the contours of vorticity for $M = 0.2, 0.4,$ and 0.6 . The periodic vortex shedding was clarified to occur. The effects of the freestream Mach number on the frequency of the vortex shedding are shown in **Figure 5**. As mentioned above, it was found that the Strouhal number becomes lower as the Mach number becomes higher. The Strouhal number for $M = 0.2$

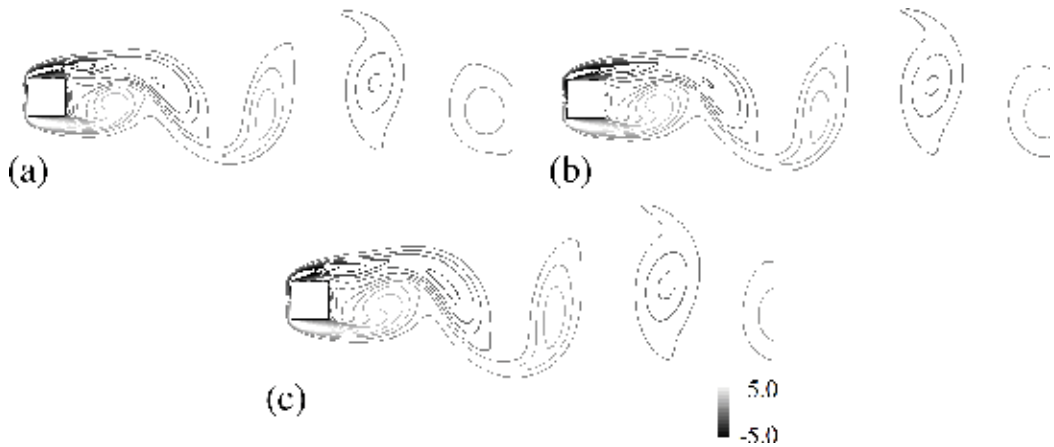


Figure 7. Contours of vorticity $\omega_z/(U_0/d)$. (a) $M = 0.2$, (b) $M = 0.4$, and (c) $M = 0.6$.

is 0.151 and that for $M = 0.6$ is 0.144. Here, to clarify the reason, the Strouhal number becomes lower, and the flow fields are discussed.

Figure 8(a) shows the mean streamwise velocity at $x/D = 1.0$ for $M = 0.2, 0.4$, and 0.6 . **Figure 8(b)** shows the half-value width of that profile, $d_{1/2}$, for $M = 0.2-0.6$. The half-value width is shown to increase as the freestream Mach number becomes higher. Also, **Figure 9** shows the mean streamwise Reynolds stress u_{1rms}/U_0 . This figure shows that the Reynolds stress becomes larger as the freestream Mach number becomes higher. This means that the velocity fluctuations of the vortices intensify. Due to this intensification, the recovery of the mean streamwise in the wake becomes more rapid and the wake becomes wider as mentioned above. This change is different

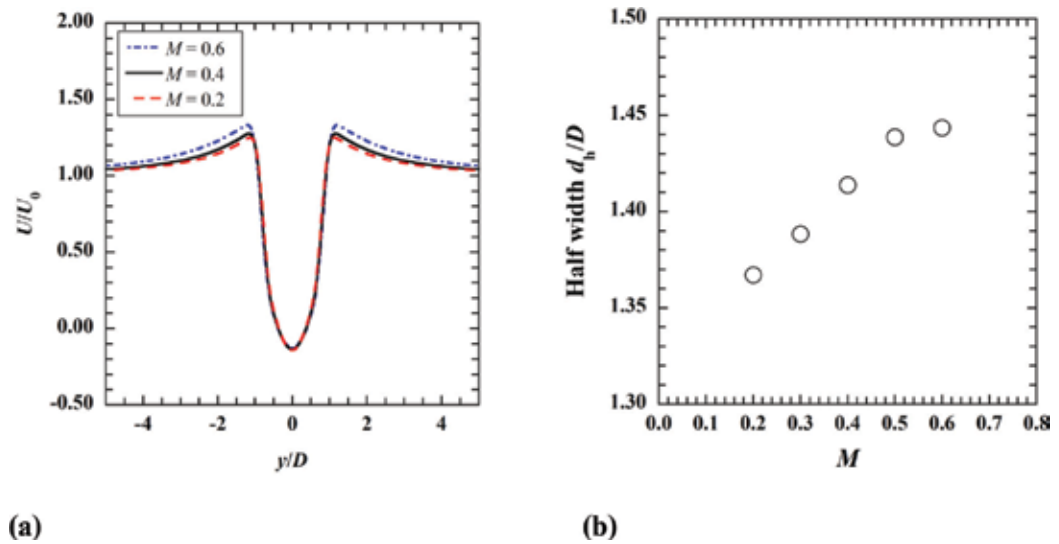


Figure 8. (a) Mean streamwise velocity at $x/D = 1.0$ ($M = 0.2, 0.4$, and 0.6) and (b) half-value width of mean streamwise velocity at $x/D = 1.0$ for $M = 0.2-0.6$.

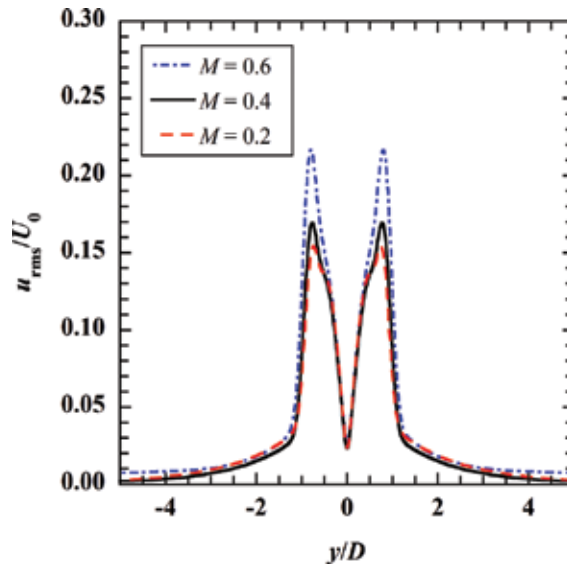


Figure 9. Mean streamwise Reynolds stress u_{1rms}/U_0 ($M = 0.2, 0.4,$ and 0.6).

from that of the vortices in the turbulent mixing layer [21], where the vortices become weaker with compressibility for a higher Mach number.

A possible reason the Reynolds stress becomes larger is that the acoustic feedback like that in the oscillations in cavity flows [22] also exists in the present cylinder flow and the acoustic waves affect the shed vortices. In this case, as the freestream Mach number becomes higher, the acoustic wave intensifies as shown in Section 4.2 and the shed vortex intensifies due to the acoustic feedback.

Roshko [23] showed that the frequency of the vortex shedding around a bluff body is proportional to the wake width. Here, to clarify the relationship between the wake width and the frequency of the vortex shedding, the modified Strouhal number St_d , which is defined by Eq. (11), was computed.

$$St_d = fd_h/U_0. \quad (11)$$

Figure 10 shows the effects of the freestream Mach number on the modified Strouhal number St_d . This figure clarifies that the modified Strouhal number is approximately independent of the Mach number. Consequently, it is confirmed that the original Strouhal number decreases because the wake becomes wider. As mentioned above, the intensification of the velocity fluctuations of the vortices widens the wake.

4.2. Acoustic radiation

Figure 11 shows the contours of pressure fluctuations with the time-averaged pressure subtracted for $M = 0.4$. For the same Mach number, Figure 12 shows the contours of the second

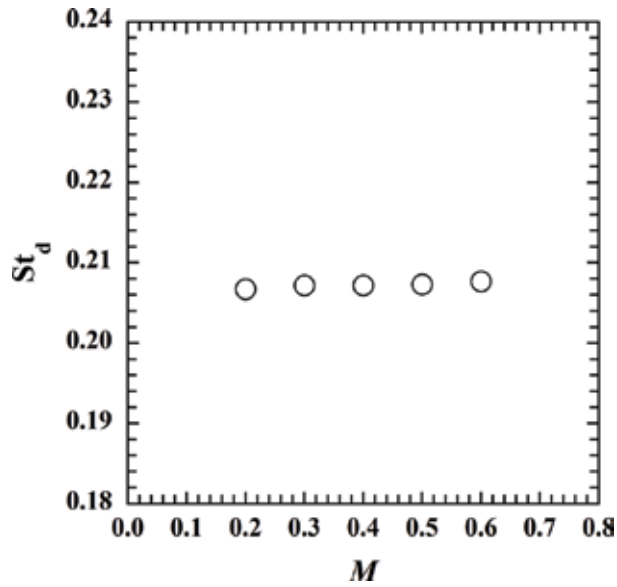


Figure 10. Modified Strouhal number.

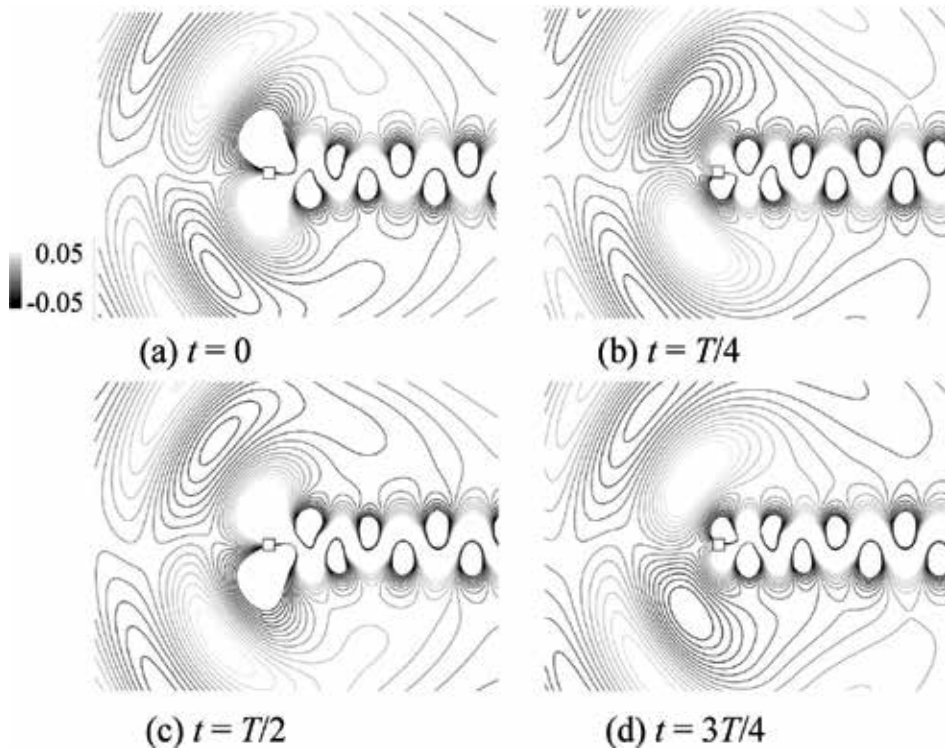


Figure 11. Contours of pressure fluctuations $p' / (\rho_0 a_0^2 M^{3.5})$ ($M = 0.4$). (a) $t = 0$, (b) $t/T = 1/4$, (c) $t/T = 1/2$, and (d) $t/T = 3/4$ (T is the period).

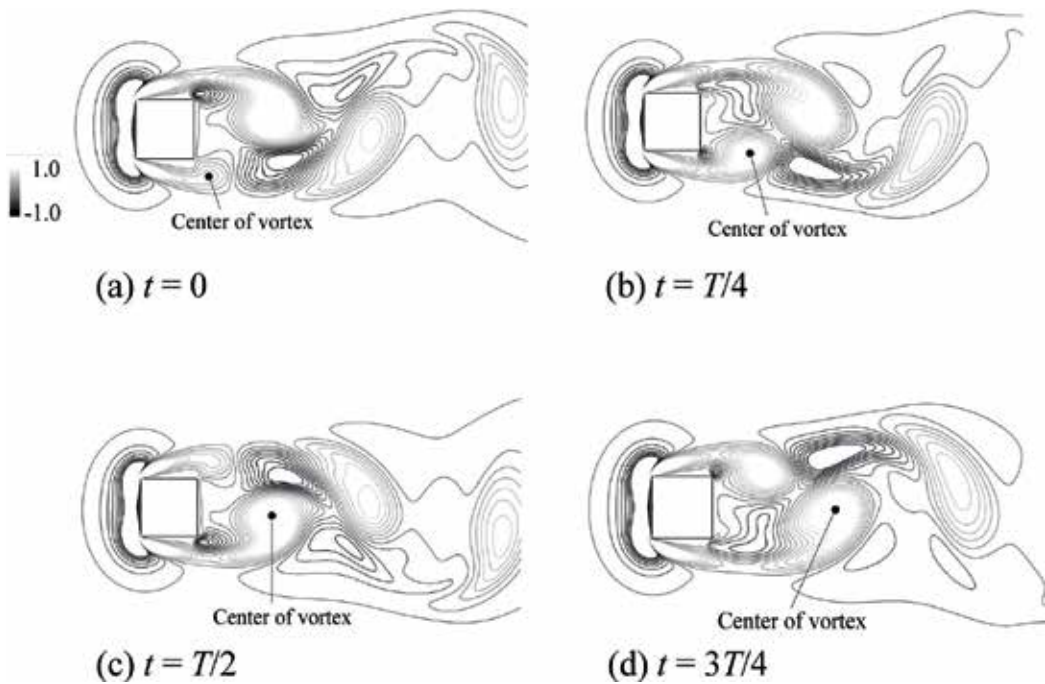


Figure 12. Contours of the second invariant of velocity gradient tensor ($M = 0.4$). (a) $t = 0$, (b) $t/T = 1/4$, (c) $t/T = 1/2$, and (d) $t/T = 3/4$ (T is the period).

invariant of velocity gradient tensor $q = \|\mathbf{\Omega}\|^2 - \|\mathbf{S}\|^2$ is computed, where $\mathbf{\Omega}$ and \mathbf{S} are, respectively, the asymmetric and symmetric parts of the velocity gradient tensor. Regions with $q > 0$ represent vortex tubes. These figures show that when the vortex is shed from the cylinder, an expansion wave is radiated on that side. For example, the expansion wave is radiated from the lower side of the cylinder in **Figures 11(b)** and **12(b)**. Meanwhile, a compression wave is radiated from the other side. This relationship of the vortex shedding and acoustic radiation is consistent with the computational results of flows around a circular cylinder by Inoue and Hatakeyama [6]. The acoustic radiation mechanism is discussed in detail.

Figure 13 shows the time histories of the pressure and density at the center of a shed vortex, where the positions of the vortex center are estimated by the local maxima of the second invariant and indicated in **Figure 12**. The pressure and density are nondimensionalized by the values at $t = 0$ (the time of $t = 0$ corresponds to **Figures 11(a)** and **12(a)**). Also, the density is raised to the power of the specific ratio γ . This figure shows that the variation of the density is approximately in good agreement with that of pressure. This means that these phenomena are adiabatic. Also, **Figure 13** shows that both the pressure and density become lower as the vortex is developed from $t = 0$ (**Figures 11(a)** and **12(a)**) to $t = T/4$ (**Figures 11(b)** and **12(b)**). This means that the fluid in the vortex expands. As a result, an expansion wave is radiated when the vortex is shed. After the shedding, the density in the vortex becomes higher and recovers to the initial value. At this time, a compression wave is radiated between the expansion waves. This radiation mechanism is independent of the freestream Mach number.

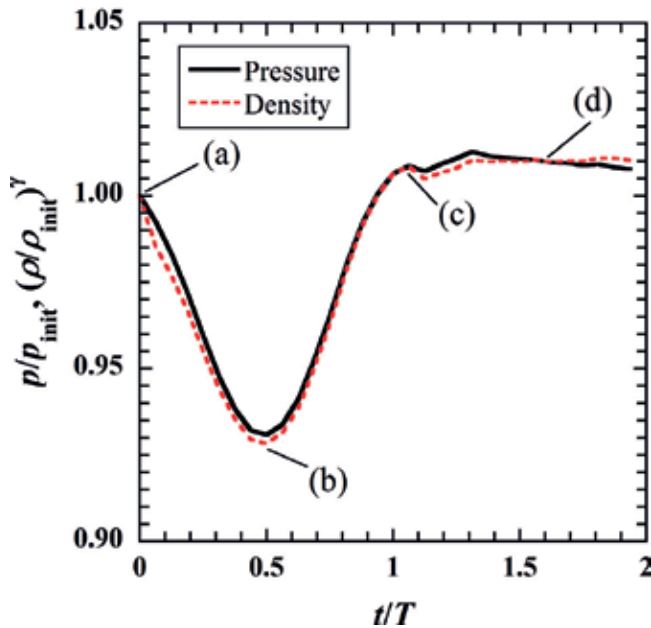


Figure 13. Time histories of pressure and density at the center of the vortex. The letters in this figure correspond to the captions of **Figures 11** and **12**.

4.3. Acoustic fields

4.3.1. Directivity of acoustic wave

Figure 14 shows the contours of the pressure fluctuations and the propagation angle of the peak of the acoustic wave, which is referred to as the propagation angle in the following. The propagation angle is compared with the theoretical angle proposed by Inoue and Hatakeyama [6]. In this theory, the scattered sound in Curle’s equation [3] is assumed to be dominant, and the sound speed is assumed to be varied by the Doppler effects as indicated in Eq. (12).

$$a_{\theta}(\theta) = a_0 \left(\sqrt{1 - M^2 \sin^2 \theta} - M \cos \theta \right), \quad (12)$$

where θ is an angle as shown in **Figure 1**. For $M = 0.2$, the propagation angle $\theta = 75^\circ$ is approximately in good agreement with the theoretical angle $\theta = 79^\circ$. However, the propagation angle $\theta = 80^\circ$ greatly differs from the theoretical angle $\theta = 62^\circ$ for $M = 0.6$. This is because the direct sound in Curle’s equation [3] becomes more intense as the freestream Mach number becomes higher. The contributions of direct and scattered sounds to total sound are presented quantitatively in Section 4.3.2.

4.3.2. Decomposition of scattered and direct sounds

The sound predicted by the direct simulation is decomposed into scattered and direct sounds. The direct sound p_{direct} is defined as.

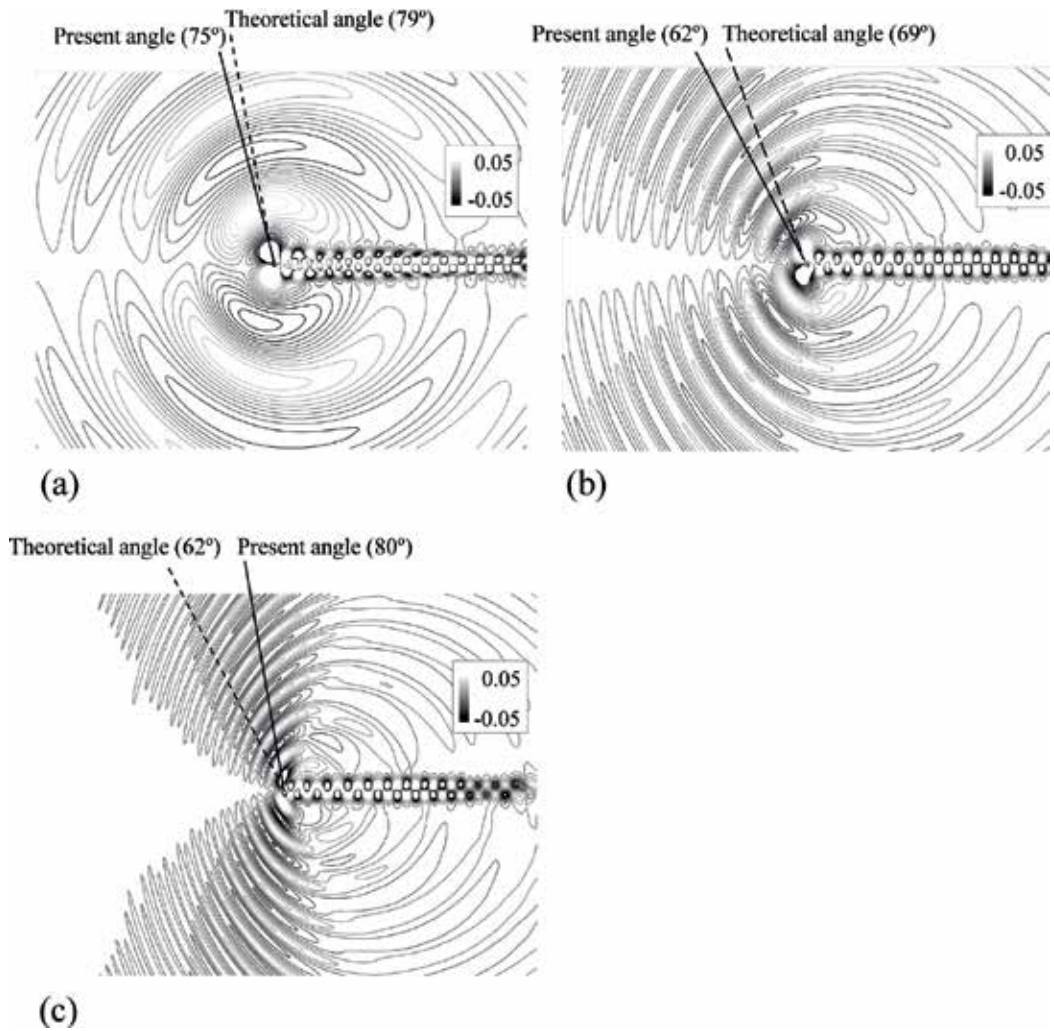


Figure 14. Contours of pressure fluctuations $p'/(\rho_0 a_0^2 M^{3.5})$ and propagation angle. Here, (a) $M = 0.2$, (b) $M = 0.4$, and (c) $M = 0.6$.

$$p_{\text{direct}}(t) = p_{\text{total}}(t) - p_{\text{scatter}}(t), \tag{13}$$

where the p_{scatter} is the dipole sound that contains the Doppler effect [6]. Also, the acoustic wavelength is $11.6D$ and so is sufficiently large even for $M = 0.6$ to neglect the difference in the retarded time on the cylinder. The scattered sound p_{scatter} is

$$p_{\text{scatter}} = \frac{1}{2^{3/2} \pi a_0^{1/2} r^{1/2}} \int_{\tau_0}^{\tau} \frac{F'(\tau')}{\sqrt{\tau - \tau'}} d\tau', \tag{14}$$

$$F' = F'_x \cos \theta + F'_y \sin \theta.$$

where $\tau = t - r/a_0$ and F'_x and F'_y are the time derivatives of the forces exerted on the fluid by the cylinder. Also, the start point of the time integration τ_0 is set to $\tau - 10 T$ to enlarge the integration interval sufficiently.

Figure 15 shows the polar plots of pressure levels of the total, scattered, and direct sounds at $r/D = 30.0$ for $M = 0.2, 0.4,$ and 0.6 . Figure 16 shows the effects of the freestream Mach number on each sound pressure levels at $r/D = 30.0$ in the above-mentioned propagation angle as shown in Table 1. For $M > 0.3$, the sound pressure level of the total sound is proportional to M^7 , although that is proportional for M^5 for $M \leq 0.3$. According to the two-dimensional Curle's

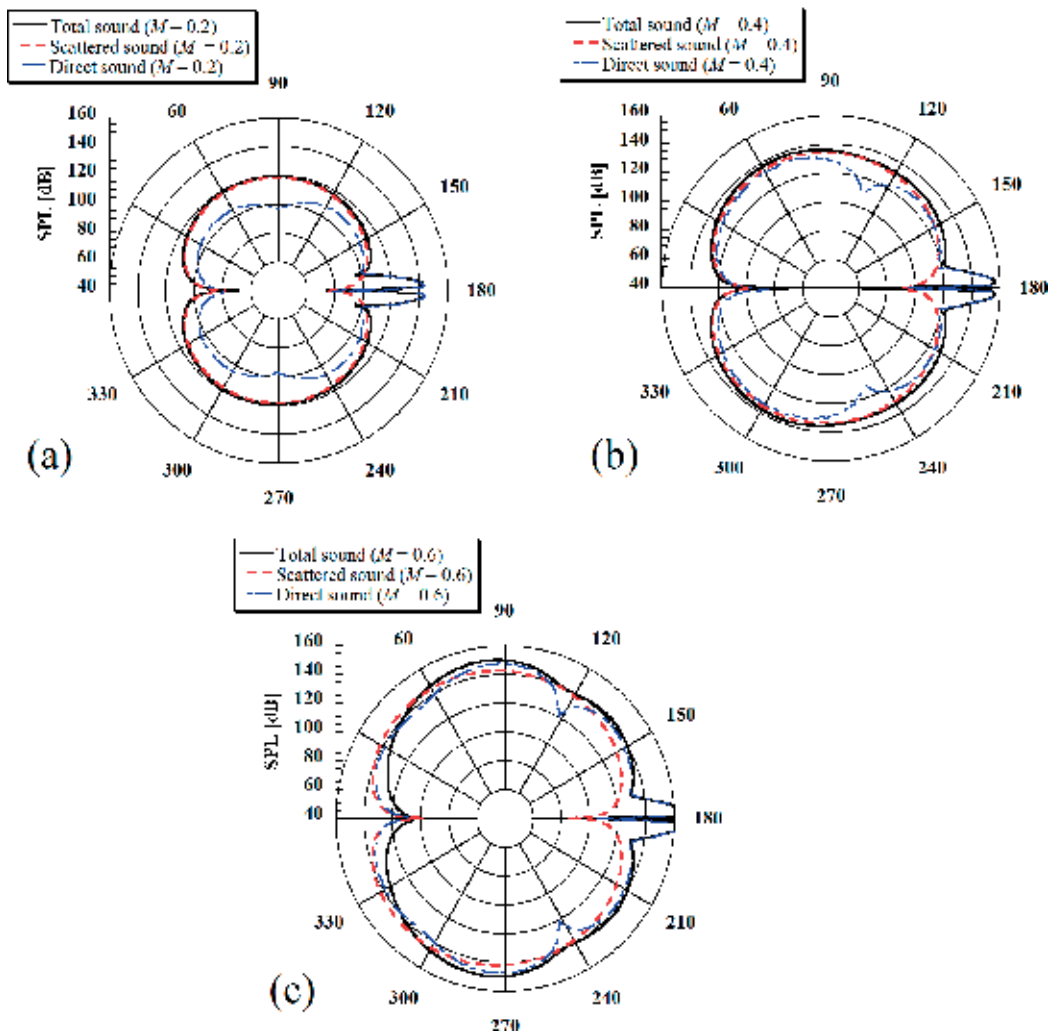


Figure 15. Polar plots of total, scattered, and direct acoustic fields at $r/D = 30.0$. (a) $M = 0.2$, (b) $M = 0.4$, and (c) $M = 0.6$ [24].

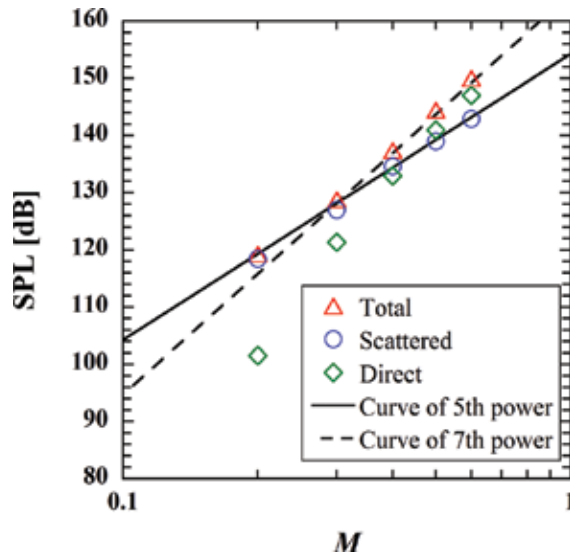


Figure 16. Effects of freestream Mach number on sound pressure levels of total, scattered, and direct sounds at the frequency of the vortex shedding at $r/D = 30.0$ in the direction of acoustic propagation angle [24].

Mach number	0.2	0.3	0.4	0.5	0.6
Present propagation angle	75	64	62	67	80

Table 1. Propagation angle of acoustic waves.

equation introduced by Inoue and Hatakeyama [6], the pressure level of the direct sound is proportional to M^7 , whereas that of the scattered sound is proportional to M^5 . The present results clarified that sound pressure levels of scattered and direct sounds intersect around $M = 0.4$.

Consequently, it was confirmed that the effects of the direct sound need to be taken into consideration when predicting the sound radiating from a cylinder flow for $M \geq 0.4$. Also, as shown in **Figure 14**, the directivity of the acoustic field for such a high Mach number cannot be predicted by the modified Curle’s equation [6], which assumes the scattered sound to be dominant and takes the Doppler effects into consideration. To the authors’ knowledge, this is the first time that the effects of the freestream Mach number on the contributions of the scattered and direct sounds have been quantitatively clarified for flows around a cylinder.

4.4. Lighthill’s acoustic sources

The right-hand term of Lighthill’s equation [Eq. (6)] can be decomposed into three components,

$$\frac{\partial^2 \tilde{T}_{ij}}{\partial x_i \partial x_j} = \frac{\partial^2 \tilde{T}_{ij}^1}{\partial x_i \partial x_j} + \frac{\partial^2 \tilde{T}_{ij}^2}{\partial x_i \partial x_j} + \frac{\partial^2 \tilde{T}_{ij}^3}{\partial x_i \partial x_j}, \quad (15)$$

$$T_{ij}^1 = \rho v_i v_j, \quad T_{ij}^2 = \delta_{ij}((p - p_0) - a_0^2(\rho - \rho_0)), \quad T_{ij}^3 = -\tau_{ij}. \quad (16)$$

Figure 17 shows the contours of the total Lighthill's acoustic sources $\partial^2 \tilde{T}_{ij} / \partial x_i \partial x_j / (\rho_0 U_0^2 / D^2)$ at the frequency of the vortex shedding in (a), those of the first term $\partial^2 \tilde{T}_{ij}^1 / \partial x_i \partial x_j / (\rho_0 U_0^2 / D^2)$ in (b), and those of the second term $\partial^2 \tilde{T}_{ij}^2 / \partial x_i \partial x_j / (\rho_0 U_0^2 / D^2)$ in (c) (hereafter referred to as first and second terms, respectively). Here, the third term is negligibly small, so its contour is not presented here. All the contours show that the acoustic sources near the cylinder are more intense than the acoustic sources in the wake far from the cylinder. This is because the acoustic waves are radiated by the vortex shedding from the cylinder as mentioned above. Also, the intensity of the second term, which is usually neglected for the acoustic prediction using Lighthill's acoustic analogy [7, 8], was found to be comparable to that of the first term.

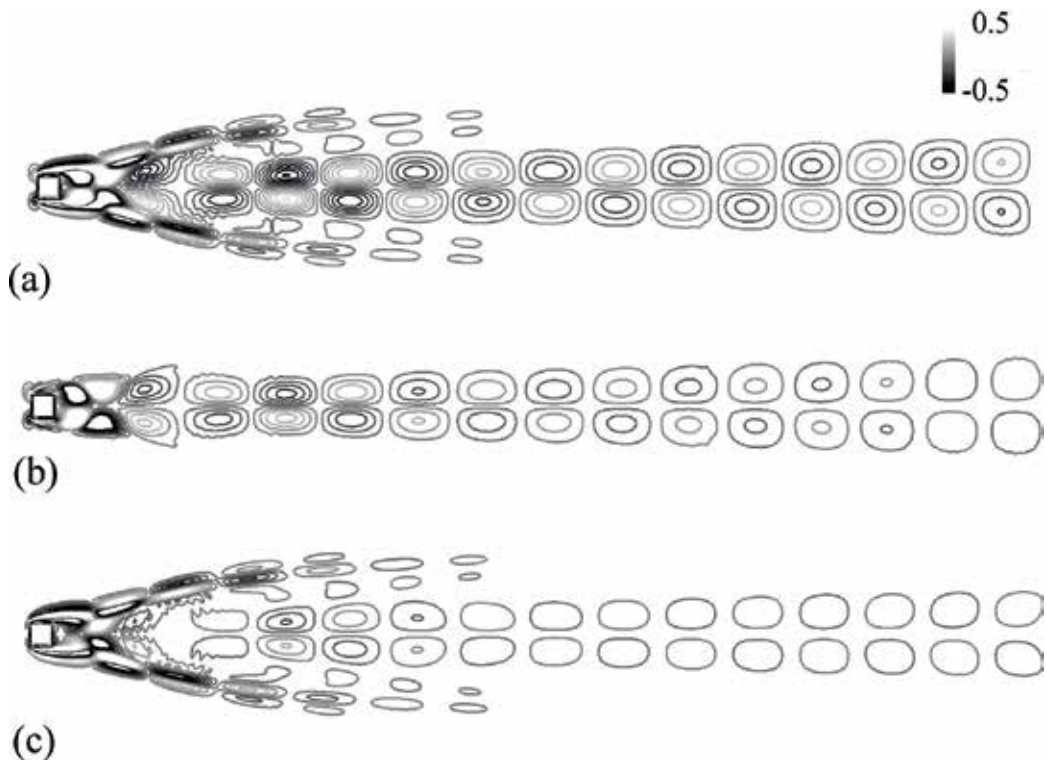


Figure 17. Lighthill's acoustic sources for $M = 0.4$ (real part is shown). (a) All terms $\partial^2 \tilde{T}_{ij} / \partial x_i \partial x_j / (\rho_0 U_0^2 / D^2)$, (b) first term $\partial^2 \tilde{T}_{ij}^1 / \partial x_i \partial x_j / (\rho_0 U_0^2 / D^2)$, and (c) second term $\partial^2 \tilde{T}_{ij}^2 / \partial x_i \partial x_j / (\rho_0 U_0^2 / D^2)$.

To clarify the contributions of each term to the acoustic field, four hybrid simulations were performed for each Mach number on the basis of total Lighthill's acoustic sources computed by the direct simulation, only the first term, only the second term, and only the third term.

Figure 18 shows the polar plots of the sound pressure level at the frequency of the vortex shedding predicted at $r/D = 30.0$ by the hybrid simulations for $M = 0.4$. It was clarified that the sound pressure levels predicted by the hybrid simulation based on all terms agree well with those based on only the first term. The sound pressure level based on the second term and that on the third term is negligibly weaker than that based on the first term. Meanwhile, the intensity of the second term is in itself comparable to that of the first term as mentioned above. This indicates that the radiation efficiency of the second term is weaker than that of the first term.

Figure 19 shows the predicted sound pressure level at $r/D = 30.0$ in the direction of the above-mentioned acoustic propagation angle. The results clarified that the first term is the dominant acoustic source for all the Mach numbers. The difference between the sound pressure level based on the first term and that based on the second or third term was more than 30 dB for all the Mach numbers. This result shows that the momentum (the first term) of Lighthill's acoustic source is the dominant acoustic source for all the Mach numbers for cylinder flows, while it has

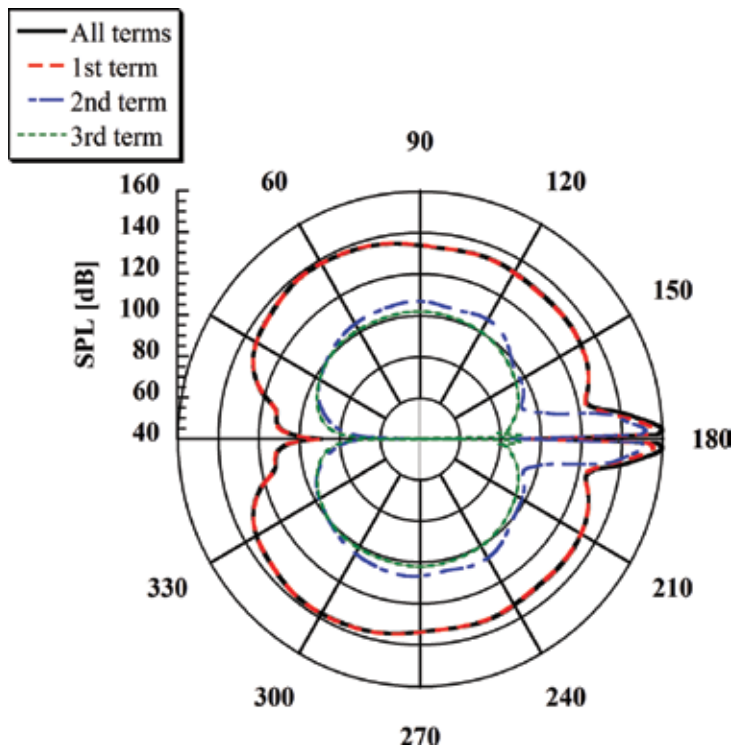


Figure 18. Polar plots of sound pressure levels predicted by decoupled simulations based on all, first, second, and third terms of Lighthill's acoustic sources at the frequency of vortex shedding at $r/D = 30.0$ [24].

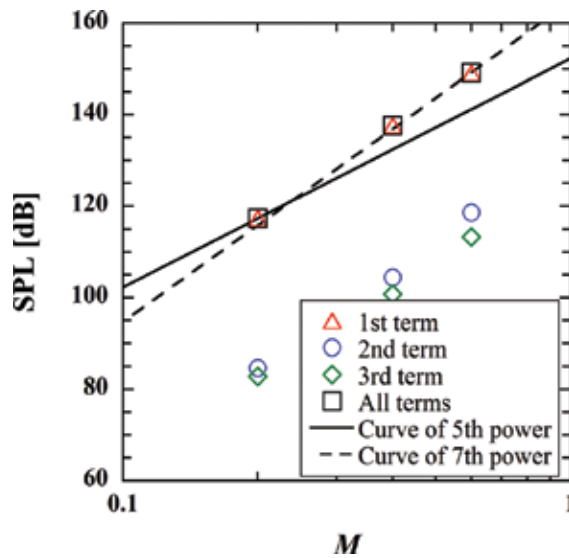


Figure 19. Effects of Mach number on sound pressure levels predicted by hybrid simulations based on all, first, second, and third terms of Lighthill’s acoustic sources at the frequency of vortex shedding at $r/D = 30.0$ in the direction of acoustic propagation angle [24].

been clarified in the past research that the entropy (the second term) also needs to be taken into consideration for high-speed jets such as $M = 0.9$. Consequently, it was confirmed that only the first term needs to be taken into consideration independently of the freestream Mach number when the sound radiating from a cylinder flow is predicted on the basis of the Lighthill’s acoustic analogy.

5. Conclusion

Aeroacoustic simulations composed of hybrid and direct simulations were introduced. The effects of the freestream Mach number on the flow and acoustic fields around a square cylinder were investigated. The Mach number was varied from 0.2 to 0.6. The Reynolds number based on the side length was 150. These results indicate the effectiveness and limit of the hybrid simulations.

It was found that the Strouhal number of vortex shedding, which is based on the side length, becomes lower as the freestream Mach number becomes higher. The Strouhal number for $M = 0.2$ is 0.151 and that for $M = 0.6$ is 0.144. As the Mach number increases, the velocity fluctuations of the vortices shed from the cylinder intensifies and the wake widens. The possible reason the velocity fluctuations of the vortices intensify is that the acoustic feedback exists like that in the oscillations in cavity flows. These effects can be found by the direct simulations.

The sound pressure level at the frequency of the vortex shedding in the direction of the acoustic propagation angle is proportional to M^7 for $M > 0.3$, while that is proportional to M^5

for $M \leq 0.3$. The decomposition of scattered and direct sound showed that the direct sound is too intense to neglect for $M \geq 0.4$. This indicates that the direct sound needs to be taken into consideration when predicting the flow-induced sound around the cylinder for $M \geq 0.4$. Also, the directivity of the acoustic field cannot be the above-mentioned modified Curle's equation for such a high Mach number.

Moreover, to clarify the contributions of each term of Lighthill's acoustic source to the acoustic field, acoustic simulations were performed using Lighthill's acoustic sources computed by the direct simulations. As a result, the momentum (the first term) of Lighthill's acoustic source was found to be dominant for all the Mach numbers while it has been clarified in the past research that the entropy (the second term) also needs to be taken into consideration for high-speed jets such as $M = 0.9$. Also, it was confirmed that only the first term needs to be taken into consideration independently of the freestream Mach number when the sound radiating from a cylinder flow is predicted on the basis of the Lighthill's acoustic analogy.

The present study has provided useful guidelines for predicting the aerodynamic sound on the basis of the Lighthill's acoustic analogy.

Acknowledgements

This work was supported by JSPS KAKENHI Grant Numbers JP26820044, 17 K06153 and through the application development for Post K computer (FLAGSHIP 2020) by the Ministry of Education, Culture, Sports, Science, and Technology of Japan (MEXT).

Author details

Hiroshi Yokoyama* and Akiyoshi Iida

*Address all correspondence to: h-yokoyama@me.tut.ac.jp

Department of Mechanical Engineering, Toyohashi University of Technology, Toyohashi, Aichi, Japan

References

- [1] Strouhal V. On one particular way of tone generation. *Annalen der Physik und Chemie* (Third series). 1878;5:216-251
- [2] Lighthill MJ. On sound generated aerodynamically Part I. General theory. *Proceedings of the Royal Society of London A*. 1952;211:564-587. DOI: 10.1098/rspa.1952.0060
- [3] Curle N. The influence of solid boundaries upon aerodynamic sound. *Proceedings of Royal Society of London A*. 1955;231:505-514. DOI: 10.1098/rspa.1955.0191

- [4] Gerrard JH. Measurements of the sound from circular cylinders in an air stream. *Proceedings of the Physical Society London, Section B*. 1955;**68**:453-461. DOI: 10.1088/0370-1301/68/7/307
- [5] Phillips OM. The intensity of aeolian tones. *Journal of Fluid Mechanics*. 1956;**1**:607-624. DOI: 10.1017/S0022112056000408
- [6] Inoue O, Hatakeyama N. Sound generation by a two-dimensional circular cylinder in a uniform flow. *Journal of Fluid Mechanics*. 2002;**471**:285-314. DOI: 10.1017/S0022112002002124
- [7] Gloerfelt X, Perot F, Bailly C, Juve D. Flow-induced cylinder noise formulated as a diffraction problem for low Mach numbers. *Journal of Sound and Vibration*. 2005;**287**:129-151. DOI: 10.1016/j.jsv.2004.10.047
- [8] Liow YSK, Tan BT, Thompson MC, Hourigan K. Sound generated in laminar flow past a two-dimensional rectangular cylinder. *Journal of Sound and Vibration*. 2006;**295**:407-427. DOI: 10.1016/j.jsv.2006.01.014
- [9] Howe MS. *Theory of Vortex Sound*. Cambridge: Cambridge University Press; 2003. 216 p. DOI: 10.1017/CBO9780511755491
- [10] Powell A. Theory of vortex sound. *Journal of the Acoustical Society of America*. 1967;**36**:177-195. DOI: 10.1121/1.1918931
- [11] Bodony DJ, Lele SK. Low-frequency sound sources in high-speed turbulent jets. *Journal of Fluid Mechanics*. 2008;**617**:231-253. DOI: 10.1017/S0022112008004096
- [12] Okajima A. Strouhal numbers of rectangular cylinders. *Journal of Fluid Mechanics*. 1982;**123**:379-398. DOI: 10.1017/S0022112082003115
- [13] Sheard GJ, Fitzgerald MJ, Ryan K. Cylinders with square cross-section: Wake instabilities with incidence angle variation. *Journal of Fluid Mechanics*. 2009;**630**:43-69. DOI: 10.1017/S0022112009006879
- [14] Tong XH, Luo SC, Khoo BC. Transition phenomena in the wake of an inclined square cylinder. *Journal of Fluids Structures*. 2008;**24**:994-1005. DOI: 10.1016/j.jfluidstructs.2006.08.012
- [15] Lele SK. Compact finite difference schemes with spectral-like resolution. *Journal of Computational Physics*. 1992;**103**:16-42. DOI: 10.1016/0021-9991(92)90324-R
- [16] Gaitonde DV, Visbal MR. Pade-type higher-order boundary filters for the Navier-Stokes equations. *AIAA Journal*. 2000;**38**:2103-2112. DOI: 10.2514/2.872
- [17] Thompson KW. Time dependent boundary conditions for hyperbolic systems. *Journal of Computational Physics*. 1987;**68**:1-24. DOI: 10.1016/0021-9991(90)90152-Q
- [18] Poinot TJ, Lele SK. Boundary conditions for direct simulations of compressible viscous flows. *Journal of Computational Physics*. 1992;**101**:104-129. DOI: 10.1016/0021-9991(92)90046-2

- [19] Kim JW, Lee DJ. Generalized characteristic boundary conditions for computational aeroacoustics. *AIAA Journal*. 2000;**38**:2040-2049. DOI: 10.2514/2.891
- [20] Colonius T, Lele SK, Moin P. Sound generation in a mixing layer. *Journal of Fluid Mechanics*. 1997;**330**:375-409. DOI: 10.1017/S0022112096003928
- [21] Terakado D, Nonomura T, Oyama A, Fujii K. Mach number dependence on sound sources in high mach number turbulent mixing layer. *Proceedings of 22nd AIAA/CEAS Aeroacoustics Conference*; 30 May–1 June 2016; Lyon, France: AIAA; 2016. p. 1-10
- [22] Yokoyama H, Kato C. Fluid-acoustic interactions in self-sustained oscillations in turbulent cavity flows. I. Fluid-dynamic oscillations. *Physics of Fluids*. 2009;**21**:105103-13-1105103-1. DOI: 10.1063/1.3253326
- [23] Roshko A. On the drag and shedding frequency of two-dimensional bluff bodies. *National Advisory Committee for Aeronautics Technical Note*. 1954;**3169**:1-29
- [24] Yokoyama H, Iida A. Identification of dominant acoustic sources in flows around square cylinder in a uniform flow. In: *Computational Modeling, Simulation and Applied Mathematics*. 2016; Chapter of ECTE2016:125-129

Analytical and Mathematical Analysis of the Vibration of Structural Systems Considering Geometric Stiffness and Viscoelasticity

Alexandre de M. Wahrhaftig,
Reyolando M. L. R. F. Brasil and
Lázaro S. M. S. C. Nascimento

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.75615>

Abstract

For a complete analysis of vibration, the stiffness of a structure must have two characteristics: one corresponding to conventional stiffness and the other to the geometric stiffness. Thus, the total stiffness takes form where the model to be used to represent any behavior of the material is introduced to the first part via the modulus of elasticity. The second is the geometric stiffness, through which it is possible to linearize a geometric nonlinear problem. To consider both aspects, a mathematical model based on the Rayleigh method has been elaborated. Two systems were numerically studied. First, the occurrence of resonance in the vibration of a prestressed reinforced concrete beam has been investigated. The results indicated resonant and non-resonant schemes between the natural frequency of the beam and the frequency of the engine. To the second system, the first natural frequency of a slender, 40-m-high concrete mobile phone mast, was calculated, and an evaluation of the structural collapse was performed. To the both systems, the cross section of reinforced concrete was treated by the theory for the homogenized section in order to consider the presence of the steel, and the viscoelasticity of the concrete was taken into account through a three-parameter rheological model.

Keywords: analytical mathematical analysis, numerical simulation, viscoelasticity, vibration, rheology, Rayleigh method, geometric stiffness, buckling load

1. Introduction

The dynamic characteristics of a structure depend, basically, on its stiffness and mass. With these two elements, the natural frequencies and modes of vibration of the system are determined. However, the initial stiffness of a structure can be affected by the so-called geometric stiffness, a function of the acting normal force. In the case of compression force, the stiffness of the structure decreases, also reducing the natural frequencies of vibration. A class of structures of socio-economic-strategic importance for the national industry are machine bases, which are subject to vibrations induced by the supported equipment. These vibrations can affect the safety of the structure itself and generate detrimental effects on the equipment and the quality of the manufactured product. They can also make the working ambience unsuitable for operators. All industrial sectors are subject to these problems, including oil exploration, production, and refining, mining, wind energy, atomic energy, as well as bridges and viaducts for road and rail use.

Although equipment support structures are, as a general rule, over-dimensioned, and therefore not subject to the effects of geometric stiffness, the tendency of modern structural engineering is towards increasingly slender elements, made possible by materials that are more efficient and lightweight, and having more and more powerful structural analysis capabilities. One of these features is prestressed concrete, represented by the presence of a steel bar or cable inside the structure that compresses it, the purpose of which is to reduce the effects of tension on flexion. In the case of beams subjected to periodic excitation, it is assumed that the original design has taken care to distance the natural frequencies of the system from those of the excitation, considering that, by hypothesis, the prestressing force decreases the stiffness of the element and, consequently, its natural frequencies, which may lead to unexpected, potentially dangerous resonance regimes. In the opposite direction, the presence of the prestressing can provide a form of control of this same vibration, where a resource is available to remove the structure of the resonant regime, if perceived in the preliminary stages of design. In one way or another, a satisfactory analysis solution to most engineering problems comes from a consideration that is easily implemented in analytical and numerical-computational formulations: the geometric stiffness. The influence of geometric stiffness has been studied in several contexts, both in laboratory tests and in comparison with the finite element method (FEM) [1–4].

The problem is aggravated when the material itself changes its elastic properties, such as in the case of viscoelasticity, which represents the gradual increase of deformation with time. This is a typical phenomenon of concrete structures because it is a viscoelastic material. It must be considered when verifying the stability of slender pieces compressed under the ultimate limit state (ULS), since these have their stiffness modified in function of the rheology of the material itself. It is important to consider the viscoelastic behavior of concrete structures relative to the characteristics of the structural element under study; this is necessary when verifying the stability of compressed slender pieces, since their stiffness is modified according to the rheology of the material. For this reason, in the specific case of columns in that loading condition, a premature analysis can produce undesired consequences, and the system may even collapse.

Typically, viscoelasticity representation is based on rheological models that are associated with deformations that are deferred over time. These models can be included in a static or dynamic analysis of the structures by relating them to the modulus of elasticity of the material. In the case of dynamic analysis, the stiffness of the structure must be composed of two terms, one of which corresponds to the portion of conventional stiffness and the other to the geometric stiffness parcel [5]. Thus, it is possible to introduce into the first one a modulus of elasticity that is variable over time, according to the rheological model adopted, keeping the stress level constant, and, in the second, to consider the normal stress acting on the system, which includes the self-weight of the structural element. An approximate and satisfactory solution can be found by considering viscoelasticity through flexural bending over time.

To evaluate these aspects, a numerical simulation has been performed, assuming an idealized section of a beam as an engine base. A rheological model of the three parameters has been used to obtain the variable modulus of elasticity. A model, including geometric stiffness, distributed, and concentrated masses, is derived based on the Rayleigh method and solved for a range of axial compression load values. The results made it allowed us to verify the resonant and non-resonant response of the system. A second analysis has been performed to simulate numerically the variation of the first natural frequency of vibration of an actual structure of reinforced concrete, axially loaded, and considering also the viscoelasticity by means of the same rheological model. The loss of stability of the system has been then evaluated.

2. Mathematical solution for representing the viscoelasticity

An increase in strain over time under constant stress is a viscoelastic phenomenon. Mathematically, viscoelasticity can be represented by a time-dependent function associated with rheological models capable of describing the phenomenon [6]. The slow deformation for concrete parts is a phenomenon that is related to loads and deformations but is partially reversible [7]. It is a phenomenon that is directly related to the movement of moisture inside concrete. When a sample of concrete is loaded for 90 days and then unloaded, the immediate or elastic recovery is approximately the same magnitude as the elastic deformation when the first load is applied [8].

It is conceptually convenient to consider classic viscoelastic models in which only two types of parameters, relating to elasticity and viscosity, appear [9]. Classic viscoelastic models are obtained by arranging springs and dampers, or dashpots, in different configurations. Springs are characterized by elastic moduli and dashpots by viscosity coefficients. The best known of these mechanical models are the Maxwell model, containing a spring in series with a dashpot, and the Kelvin-Voigt model, containing a spring and dashpot in parallel. One model used to represent the viscoelasticity of solids is the three-parameter model, in which the elastic parameter E_0 is connected to the viscoelastic Kelvin-Voigt model with parameters E_1 and η_1 , which is a simplification of the Group I Burgers model, as shown in **Figure 1**.

The three-parameter model sufficiently describes the viscoelastic nature of many solids and is often used to study the phenomenon in various scientific fields. The total deformations of the

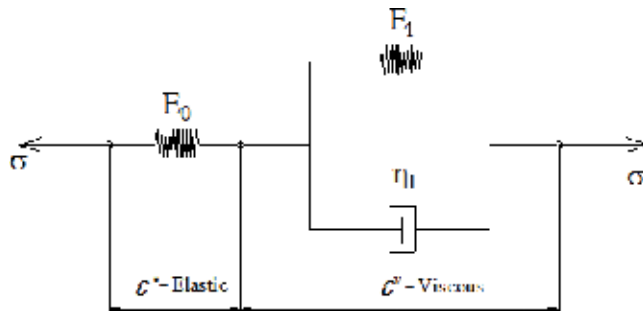


Figure 1. Viscoelastic model of three parameters.

Kelvin-Voigt model are given by $\varepsilon = \varepsilon^e + \varepsilon^v$, where ε^e is the deformation of the elastic model, and ε^v is the deformation of the Kelvin-Voigt model. When differentiated with respect to time, the total deformation is obtained as

$$\varepsilon = \dot{\varepsilon}^e + \dot{\varepsilon}^v \tag{1}$$

which is the constitutive equation of the elastic and Kelvin-Voigt models, respectively. Considering $E_1 = E_0$ as the modulus of elasticity for both parts of the rheological model,

$$\sigma = E_0 \varepsilon^e \text{ and } \dot{\sigma} + \frac{E_0 + E_0}{\eta_1} \sigma = E_0 \dot{\varepsilon} + \frac{E_0 E_0}{\eta_1} \varepsilon \tag{2}$$

are found. From the previous equations, one derives the following differential equation:

$$\sigma = E_0 \varepsilon^v + \eta_1 \dot{\varepsilon}^v \tag{3}$$

where $\sigma = 0$ for $t < 0$ and $\sigma = \sigma_0$ for $t > 0$, with t representing the time and $t = 0$ the instant of loading application. As the stress remains constant, the stress derivate with respect to time is zero. Applying the previous stress condition, the following ordinary differential equation is found:

$$E_0 \dot{\varepsilon} + \frac{E_0 E_0}{\eta_1} \varepsilon = \sigma_0 \tag{4}$$

for which the general solution for $t > 0$, taking the initial condition $\varepsilon(0) = \sigma_0/E_0$, is

$$\varepsilon(t) = \sigma_0 \left[\frac{1}{E_0} + \frac{1}{E_0} \left(1 - e^{-\frac{E_0}{\eta_1} t} \right) \right] \tag{5}$$

Obviously, if the stress level remains constant, the modulus of elasticity should decrease concurrently with increasing strain:

$$E(t) = \frac{1}{\frac{1}{E_0} + \frac{1}{E_0} \left(1 - e^{-\frac{E_0}{\eta_1} t} \right)} \tag{6}$$

The previous solution for consideration of the viscoelastic behavior of materials was used by [10] to evaluate the stability of a slender wooden column, for example. However, it is of interest, at this moment, to make clear that the present work is a numerical approximation, which takes into account the viscoelastic behavior of the concrete by assuming a viscoelastic rheological model of three parameters or as also is known of the solid standard.

It is important to note that the viscoelastic behavior of the considered material is completely represented by the temporal modulus of elasticity. Therefore, the solid such behavior is wished to study should be according that adopted model. Any material can be represented by it, being, however, its usage conditioned by performing of experimental studies in order to confirm if it is correct or not. Keeping this in mind, the concrete viscoelastic behavior is assumed to be represented by the solid standard model, as an approximation of the reality. However, criteria from regulatory codes can be used in substituting of that model or even any other rheological models can be adopted.

3. Beam as a basis of supporting

3.1. Basic considerations prestressing in reinforced concrete

A piece can be considered as prestressed reinforced concrete when it is subjected to the action of the so-called prestressing forces and of permanent and variable loads, so that the concrete is not subjected to tension or it occurs below the limit of its resistance. As an example, take the normal stresses beam diagrams of the prestressed beam of **Figure 2**, where P is the prestressing force, M_p the bending moment due to eccentricity of the load P , M_p is the bending moment due to uniformly distributed load p and R is the resultant, each one of these with their corresponding normal stresses. Under the conditions presented, the lower fibers of the beam, under positive bending moment, will have the tension stresses overturned by the superposition of those produced by the normal stress of the applied stress eccentrically.

Prestressed concrete was developed scientifically from the beginning of the last century. Prestressing can be defined as the artifice of introducing, in a structure or a part, a previous state of stresses, in order to improve its resistance or its behavior in service, under the action of

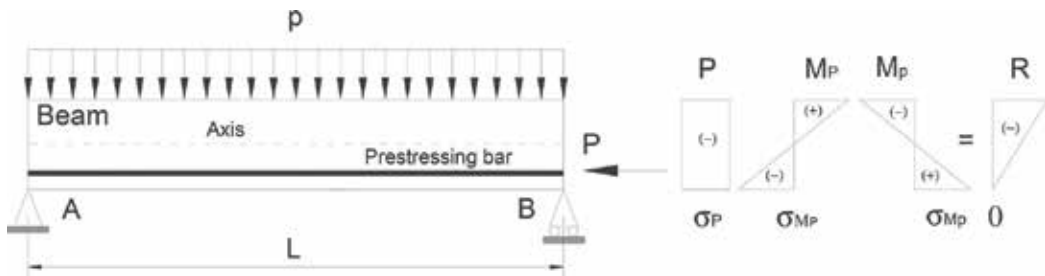


Figure 2. Normal stresses in a prestressed beam.

several effects. Due to the characteristics of the concrete as a structural material, the use of prestressing can bring a great advantage from the economic point of view. When comparing the cost of a prestressed structure with a similar one of conventional reinforced concrete, there is a reduction in the final cost of the structure due to the reduction of steel reinforcement [11]. In addition, the prestressing allows the part to overcome large spans, improves the control and reduction of deformations and fissures. It can also be used for structural recovery and reinforcement, as well as for slender systems and prefabricated or precast parts. There are three types of prestressing systems: (1) prestressing with initial adherence, (2) prestressing with posterior adherence, and (3) prestressing without adherence. The latter type is composed of a post-tensioning system characterized by the slipping freedom of the steel reinforcement in relation to the concrete, along the whole extension of the cable, except for the anchorages.

In a non-adherent prestressing, the cables or chutes are wrapped in two or three layers of resistant paper. The wires and paper are painted with bituminous paint in order to tension them after the concrete has hardened. The bitumen avoids the penetration of the cement cream inside the cable and, in this way, it eliminates adhesion between the concrete and the reinforcement [12]. The prestressed concrete is a composite material of the aggregate mixture and a cement paste associated with prestressing cables and/or passive reinforcing bars. Because of the combination of several materials, these structures develop a highly complex behavior, presenting a non-linear response, which is due, among other factors, to time-dependent effects, such as the creep of the concrete [13].

3.2. Mathematical model for the nonlinear vibration problem

Consider a rotary machine mounted on a beam subjected to a pre-tensioning force, without adhesion. It is known that such forces affect the geometric stiffness and, consequently, the values of the undamped free-vibration frequencies. If the structure is designed, as is usually the case, to have frequencies farther from the machine's service speed rotation, the changes in the frequency due to geometric stiffness may lead to the appearance of potentially dangerous resonance conditions.

Take a beam model of Bernoulli-Euler applied to a simply supported beam AB of length L and inertia I , intended to function as the base of an engine E_g , composed of viscoelastic material, represented by the temporal modulus of elasticity $E(t)$ as shown in **Figure 3**. A normal force of compression P reproduces the post-tensioning force, which changes the stiffness, and consequently, the natural frequency of vibration of the structure with time. The eccentricity between

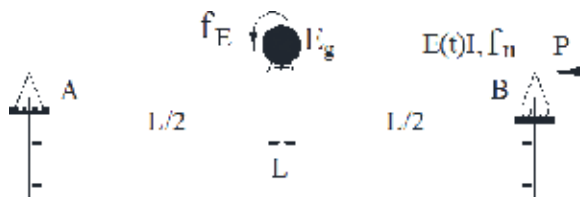


Figure 3. Beam model.

the engine axis and the part is initially ignored. The vertical displacement of the central joint is the generalized coordinate of the system.

By using the Rayleigh method [14], the undamped vibration frequency in its first mode considering viscoelasticity is obtained. It is worth mentioning that Rayleigh assumed that a system containing infinite degrees of freedom could be associated with another with a single degree of freedom (SDOF) to approximate its frequency. It is important to note that the technique developed by Rayleigh aimed at calculating the fundamental vibration frequency of elastic systems, as its precision is dependent on the function chosen to represent this mode of vibration.

The basic concept of the method is the principle of energy conservation, and can, therefore, be applied to linear and nonlinear structures. [15] applies the Rayleigh technique to determine the fundamental period of vibration to verify the stability of mechanical systems. The process is described in relation to the principle of virtual works and as the appropriate choice of the generalized coordinate describing the first mode of vibration. At the end of the process, the generalized properties of the system are obtained as stiffness and mass, necessary for the calculation of the frequency.

Consider that the vertical displacement of a generic section of the beam in **Figure 4** is given by:

$$v(x, t) = \phi(x) q(t), \tag{7}$$

in which $\phi(x)$ is a shape function that attempts to define the boundary conditions in the supports and value 1 in the central section of the beam, whose displacement with time is $q(t)$. In this case, one adopts the shape function $\phi(x) = \sin(\pi x/L)$, which is the exact solution of the problem without the P load. A prime mark will denote a derivative of the function in relation to x (Lagrange's notation).

Applying the Rayleigh method, one has the conventional bending stiffness, K_0 , as a function of the material behavior and the geometry of the cross, which is equivalent to:

$$K_0(t) = \int_0^L E(t)I(\phi'')^2 dx = \frac{\pi^4 E(t)I}{2L^3} \tag{8}$$

where $E(t)I$ is the known flexural bending with viscoelasticity, represented by multiplication of the temporal material modulus of elasticity with the inertia of the section in relation to the considered movement, the vertical vibration mode (1st mode). In turn, the geometric stiffness, K_G , as a function of the normal force of compression (or even tension), is equivalent to:

$$K_G = P \int_0^L (\phi')^2 dx = \frac{P\pi^2}{2L} \tag{9}$$

The total generalized mass of the system is found by calculating $M = M_C + M_V$ where M_C is the concentrated mass at the middle span and M_V is the mass coming from the beam self-weight given by:

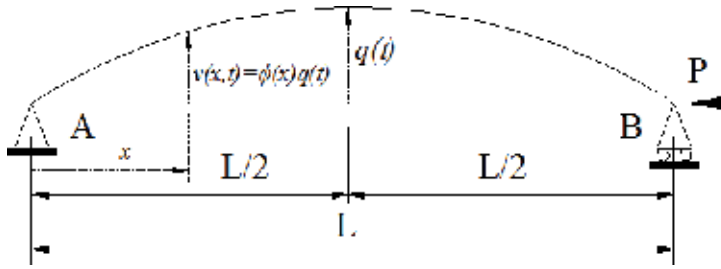


Figure 4. Rayleigh method.

$$M_v = \int_0^L m_V \phi(x)^2 dx = \frac{m_V L}{2} \tag{10}$$

in which m_V represents the total mass per length unit. Finally, the frequency of undamped free vibration (in rad/s) is found by way of Eq. (11):

$$\omega(t) = \sqrt{\frac{K(t)}{M}} \tag{11}$$

Considering the total beam stiffness as $K(t) = K_0(t) - K_G$, the free undamped frequency of vibration of the 1st mode is found, in Hertz, admitting the compressive force as positive, by:

$$f(t) = \frac{\omega(t)}{2\pi} = \frac{1}{2} \left[\frac{\pi^2 E(t) I - PL^2}{L^3 (L m_V + 2Mc)} \right]^{\frac{1}{2}} \tag{12}$$

For a better understanding of the Rayleigh method and the importance of the geometric stiffness to the structural analysis, the work [16, 17] should be consulted.

3.3 Numerical simulation 1

The beam gross cross section was estimated with a passive reinforcement arrangement capable of resisting the predicted load in the simulation, being treated by the homogenized section method, with geometry as indicated in **Figure 5**. The modulus of elasticity of the concrete was calculated according to NBR 6118/2014 [18] recommendations, following Eq. (13), for a concrete characteristic compressive strength, f_{ck} , equal to 30 MPa.

$$E_0 = \alpha_i \cdot 5600 \sqrt{f_{ck}} = 26838.405 \text{ MPa};$$

$$\alpha_i = 0.8 + 0.2 \cdot \frac{f_{ck}}{80 \text{ MPa}} = 0.875 \tag{13}$$

The reinforced concrete specific weight γ_c was obtained for a material density ρ of 2500 kg/m³ and a gravitational acceleration g of 9.8061 m/s², therefore $\gamma_c = 24.52 \text{ kN/m}^3$.

Section data:

- External height: $H = 16$ cm
- Internal height: $h = H - 2 \cdot t = 6$ cm
- Wall thickness: $t = 5$ cm
- Concrete cover: $c' = 25$ mm
- Reference beam span: $L = 3$ m
- Reinforced bar diameters: $d_b = 8$ mm
- Number of reinforcement bars: $nb = 4$
- Total inertia: $I = \frac{H^4 - h^4}{12} = 5353$ cm⁴
- Gross-sectional area: $A = H^2 - h^2 = 220$ cm²

The concrete section was homogenized by the transformation of the steel bars of the reinforcement, which led to a homogenization factor of 1.061 to be considered in the material and geometric properties of the beam section. The homogenizing technic is presented in Section 4.2. For the simulation, all the elements that constitute physical parts to be added to the system, such as the bar used in prestressing systems and an electric induction motor that represents periodic excitation, were considered as lumped or distributed masses.

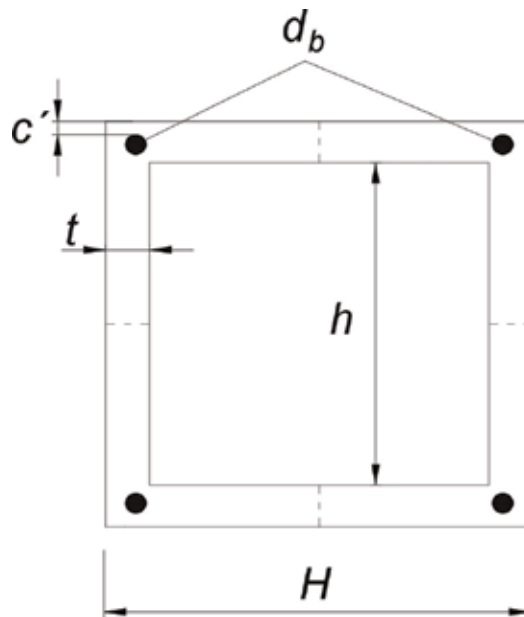


Figure 5. Beam section characteristics.

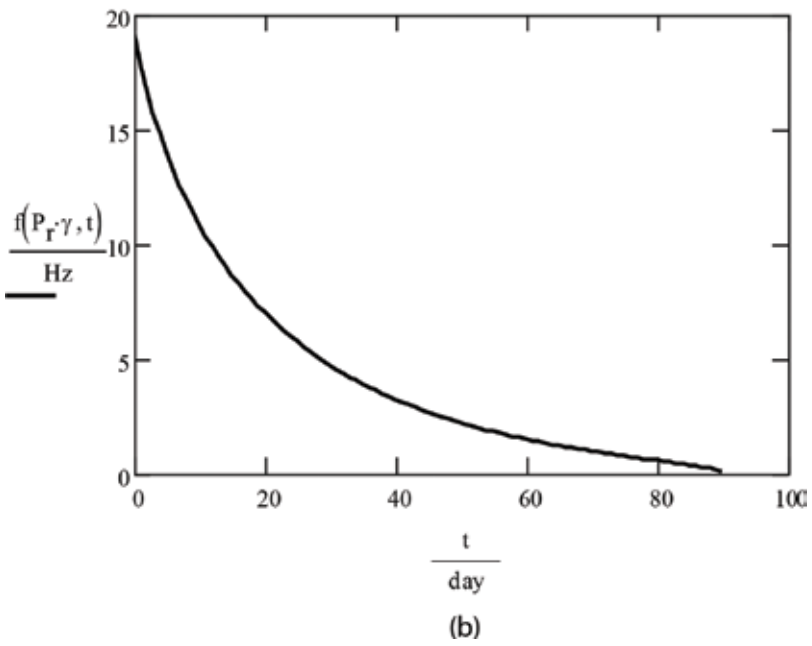
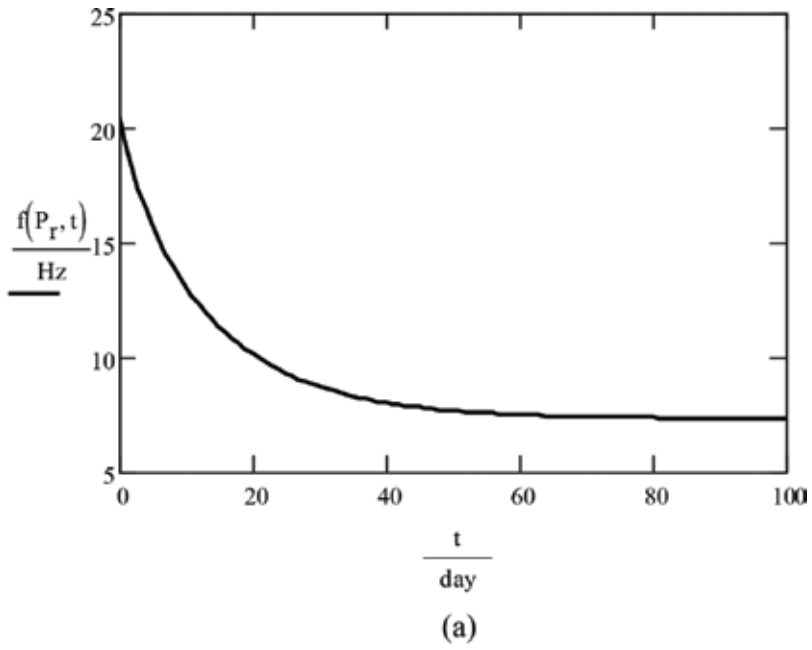


Figure 6. Frequency of the beam with time with viscoelasticity. (a) Frequency with time for capacity of the section, (b) Safety factor γ —Collapse at 90th day.

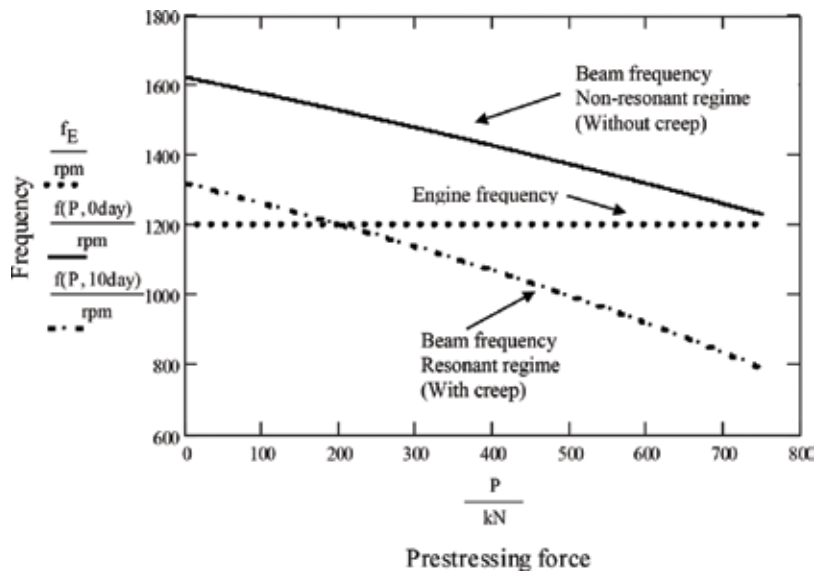


Figure 7. Resonant and non-resonant frequencies as a function of axial compression force P .

By fixing the force on the section-resistant capacity, one can observe the variation of the natural frequency of the beam with time when considering viscoelasticity, as shown in **Figure 6(a)**. A safety factor of 1.170426 to be applied to the loading can be found, which defines the beam collapse at the 90th day, as can be seen in **Figure 6(b)**. By varying force P , which represents a non-adherent post-tension force, from zero to the resistant capacity of the section, the variation of the natural frequency of the beam is obtained, given in the graph in **Figure 7**. There, it is possible to see that, with the increase of the axial compression force, the beam frequency decreases, since the geometric portion (K_G) stiffness is changed, consequently decreasing the total stiffness (K) of the beam.

Since the motor rotation is set at 1200 rpm (20 Hz), represented by the dotted horizontal line, there is no resonance without consideration of viscoelasticity, but the resonance appears when the natural frequency of the beam is calculated with the introduction of the viscoelastic behavior of the material. For 10 days after application of load, for example, the resonant regime can be observed by the intersection of two curves, dotted (horizontal) and dashed (sloped), defining exactly for which prestressing force the phenomenon occurs at that instant (210 kN).

4. Column as a structure for transmitting system

4.1. Formulation of the undamped vibration problem

Assuming the well-known trigonometric function

$$\phi(x) = 1 - \cos\left(\frac{\pi x}{2L}\right), \tag{14}$$

where x is an independent variable of the problem originating in the base, in the cantilever position, and L is length of the column, $q(t)$ is the generalized coordinate, and $e(t)$ is the vertical displacement of the top due the vibratory movement, as shown in **Figure 8**. By using the Rayleigh method, in a similar way as described in the previous Section 3.2, the conventional stiffness is found by

$$k_{0s}(t) = \int_{L_{s-1}}^{L_s} E(t)I_s F_s (\phi'')^2 dx, \text{ with } K_0(t) = \sum_{s=1}^n k_{0s}(t), \tag{15}$$

where $k_{0s}(t)$ and F_s are the parcel of the stiffness and the homogenizing factor of the concrete cross section due to the reinforcement steel at the segment s . K_0 is the final conventional stiffness, where n is the total number of intervals given by the structural geometry. $E(t)$ represents the variable modulus of elasticity on time, according Eq. (6), and I_s is the moment of inertia of each section.

The geometric stiffness is obtained by the following equation:

$$k_{gs} = \int_{L_{s-1}}^{L_s} N(x) (\phi'')^2 dx, \text{ with } K_g = \sum_{s=1}^n k_{gs}, \tag{16}$$

where k_{gs} is the geometric stiffness at the interval s ; K_g is the total geometric stiffness of the structure, with n as defined before; $N(x)$ is a normal effort function at the respective interval,

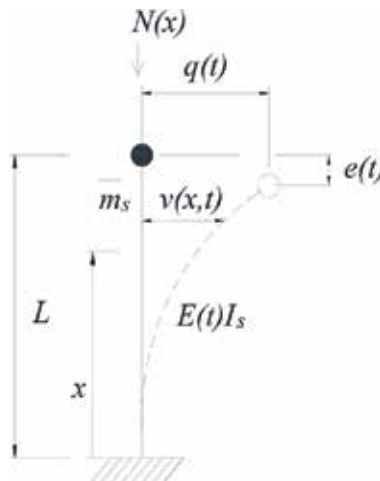


Figure 8. Mathematical model of vibration of a column.

which includes the self-weight of the column on considered part, and the lumped forces from upper segments, or better

$$N(x) = \{m_0 + \bar{m}_s[(L_s - L_{s-1}) - x]\}g, \tag{17}$$

with m_0, \bar{m}_s are defined below, and g is the acceleration of gravity. The generalized mass is then given by $M = m_0 + m$, where m_0 is the lumped mass on the top and m is the generalized mass found with

$$m_s = \int_{L_{s-1}}^{L_s} \bar{m}_s (\phi(x))^2 dx, \text{ with } \bar{m}_s = A_s \rho \text{ and } m = \sum_{s=1}^n m_s, \tag{18}$$

where \bar{m}_s is the mass to each segment s , found by multiplying the cross-sectional area A_s to the density ρ of the material at the respective intervals, that is, \bar{m}_s , mass per unit length, and m , generalized mass of the system due the density of the material, with n as defined previously. The first natural frequency of the structure is calculated:

$$\omega(t) = \sqrt{\frac{K(t)}{M}}(rd/s):f(t) = \frac{\omega(t)}{2\pi} (\text{Hertz}), \tag{19}$$

taking into consideration that, for a compressive force being positive, the temporal stiffness is:

$$K(t) = K_0(t) - K_g. \tag{20}$$

It is important to mention that Eq. (14) has been evaluated by [19] as a valid shape for the first mode of vibration with geometric nonlinear characteristics, applied for actual structures, even those with variable geometry, being a function valid throughout the entire domain of the structure.

4.2. Numerical simulation 2

A 40-m-high reinforced concrete pole structure with an external 60-cm hollow circular cross-sectional diameter, with variable thickness (**Figure 9**) and a slenderness ratio of 472 was used for analysis. The properties of the sections change along the length due to the changes in thickness and variation of the steel area.

The concrete used in the manufacture of the structure had the compression characteristic strength f_{ck} of 45 MPa, viscosity η_1 of 51089681149.92 MPa·s, and a density ρ of 2600 kg/m³. The concrete cover c' specified for the reinforcing steel was 25 mm and the steel used in the construction of the pole was CA-50, with yield strength of 500 MPa and modulus of elasticity E_s of 205 GPa. The secant modulus of elasticity of the concrete E_c is 31931.05 MPa. The numerical simulation was performed considering that all elastic parameters in Eq. (6) are equal to the modulus of elasticity of the concrete, $E_0 = E_1 = E_c$.

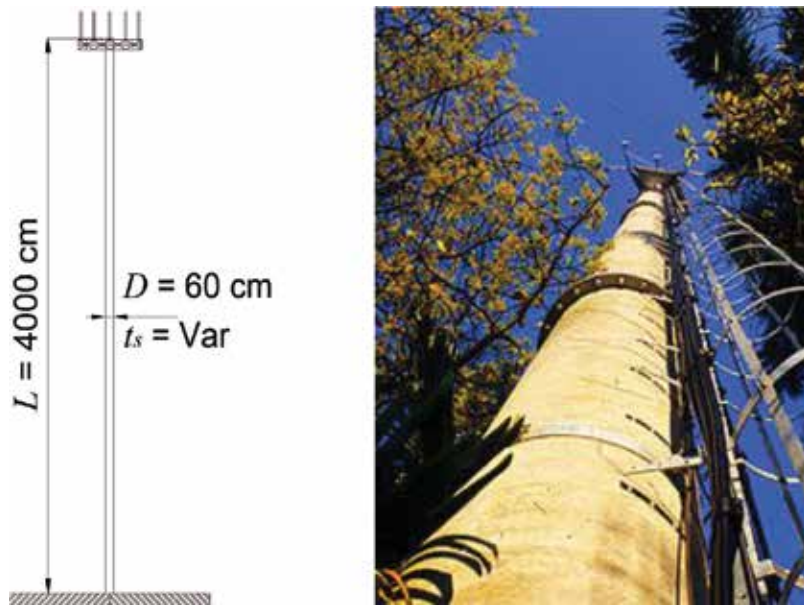


Figure 9. A column as a mast transition system.

The structure also has an array of antennas and accessories, such as a platform, stairs, cables, and mats (characteristics shown in **Table 1**), which exert compressive forces. It is important to mention that the viscous parameter was adjusted so that the deformations converged at 90 days (**Figure 10(a)**), as observed by [8]. With this, it was possible to obtain the variable modulus of elasticity $E(t)$ (**Figure 10(b)**). The gravitational (g) acceleration was assumed to be 9.806650 m/s^2 .

Since this is a cylindrical concrete reinforcement structure, it is necessary to take into account the presence of reinforcement bars at the moment of inertia of the cross-sectional area, which must be done by homogenizing the concrete area. Considering a circular ring cross section with external diameter D ; thickness of the wall t_s ; s relative to the considerate segment of the structure; a reinforcement bar b_i any occupies a position i in the section defined by R_{b_i} and θ_i ,

Dispositive	Height	Weight
Pole	from 0 to 40 m	25.48 kN/m^3
Stair	from 0 to 40 m	0.15 kN/m
Cables	from 0 to 40 m	0.25 kN/m
Platform and supports	40 m	4.90 kN
Antennas	40 m	1.88 kN

Table 1. Structure's characteristics and devices.

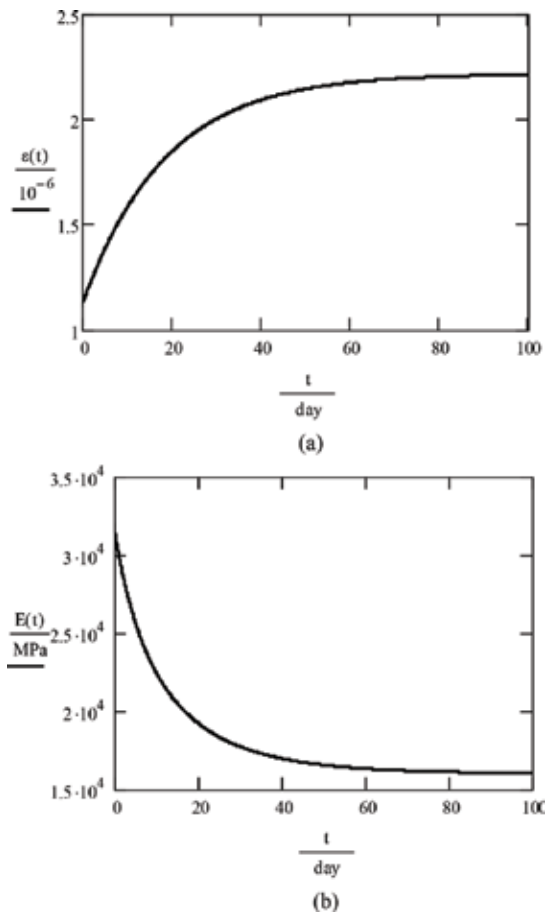


Figure 10. Deformation and modulus of elasticity over time. (a) Deformation, (b) Elasticity.

as shown in **Figure 11**. R_{bi} determines the center position of each bar in relation to the section center. c' is the concrete cover of the reinforcement and d_{bi} is the diameter of the i bar.

$$R_{bi} = \frac{D}{2} - c' - \frac{d_{bi}}{2}. \quad (21)$$

As θ_i is the independent variable, the distance between the center of each bar relative to the axis center of inertia of the section is

$$y(\theta_i) = \text{sen}(\theta_i)R_{bi}. \quad (22)$$

The spacing between the center of each bar section was obtained for $sp. = 2\pi R_{bi}/n_{bi}$, where n_{bi} is the number of bars of the reinforcement steel. The angular phase shift between them is $\Delta\theta = sp/R_{bi}$.

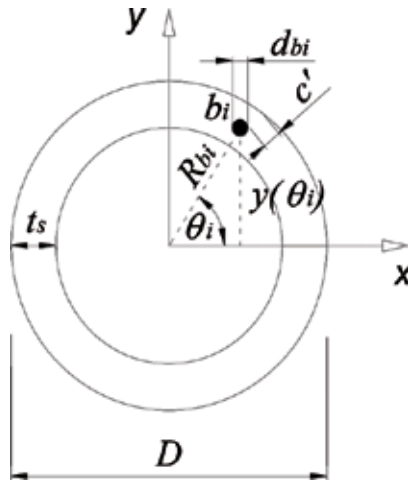


Figure 11. Parameters for homogenizing concrete section.

Since the θ_i varies from 0 to 2π at intervals defined by $\Delta\theta$, the total inertia of the steel bars in relation to the section of the structure could be obtained by the theorem of parallel axes with the expression (14).

$$I_s = \sum_{\theta}^{2\pi} \left(\frac{\pi d_{bi}^4}{64} + y(\theta_i)^2 \frac{\pi d_{bi}^2}{4} \right). \tag{23}$$

The homogenized moment of inertia of the steel bars will thus be:

$$I_{shom} = \sum_{\theta}^{2\pi} I(\theta_i) \left(\frac{E_s}{E_c} - 1 \right). \tag{24}$$

The total homogenized inertia of the section will be obtained by $I_{tot} = I + I_{shom}$, with I being the inertia of the circular section, $I = \pi/64 [D^4 - (D - 2t_s)^4]$. Thus, to find a factor F_s , which multiplies the nominal inertia of the section in terms of total steel inertia, the homogenized section is made by $F_s = 1 + (I_{shom}/I_{tot})$. Factors of homogenizing, the structural properties and the geometry of the structure are shown in **Table 2**.

Considering that the actual structure has variable proprieties along the height, the expressions (16), (17) and (18) must be resolved for each interval defined by structural geometry. The frequency was calculated for the 90th day according to Eq. (10) (see **Figure 12**).

Figure 13(a) shows the structural frequency over time, for a height limit of losing stability, calculated for the 90th day, considering viscoelasticity ($L = 50.6975$ m). To a height of 57.0000 m, for example, the behavior of **Figure 13(b)** is found, with the structural collapse occurring on the 60th day. The height limit without viscoelasticity (instant 0) is 71.29 m, a frequency of 0.0000 Hz.

Similar simulations for evaluation of the viscoelasticity can be found in [20, 21].

Height L_s (m)	External diameter D (cm)	Thickness t_s (cm)	Number of bar (n_{bi})	Bar diameter d_{bi} (mm)	Factors of homogenizing F_s
40	60	10	20	13	1.0963
39	60	10	20	13	
38	60	10	20	13	
37	60	10	20	13	
36	60	10	20	13	
35	60	10	20	13	
34	60	10	20	13	
33	60	10	20	13	
32	60	10	20	13	
31	60	13	20	13	1.0869
30	60	12	15	16	1.0995
29	60	11	15	16	1.1029
28	60	11	15	16	
27	60	11	15	16	
26	60	11	15	16	
25	60	11	16	16	1.1091
24	60	11	17	16	1.1153
23	60	11	18	16	1.1214
22	60	11	19	16	1.1274
21	60	11	20	16	1.1334
20	60	14	20	16	1.1230
19	60	15	15	20	1.1374
18	60	16	15	20	1.1354
17	60	13	16	20	1.1512
16	60	13	16	20	
15	60	13	17	20	1.1594
14	60	13	18	20	1.1675
13	60	13	19	20	1.1755
12	60	13	19	20	
11	60	13	20	20	1.1833
10	60	13	22	20	1.1987
9	60	16	22	20	1.1889
8	60	16	15	25	1.1961
7	60	17	15	25	1.194
6	60	14	16	25	1.2132
5	60	14	16	25	
4	60	14	17	25	1.2241
3	60	14	17	25	
2	60	14	17	25	
1	60	18	17	25	1.2136
0	60	18	17	25	

Table 2. Structural properties and homogenizing factors of sections.

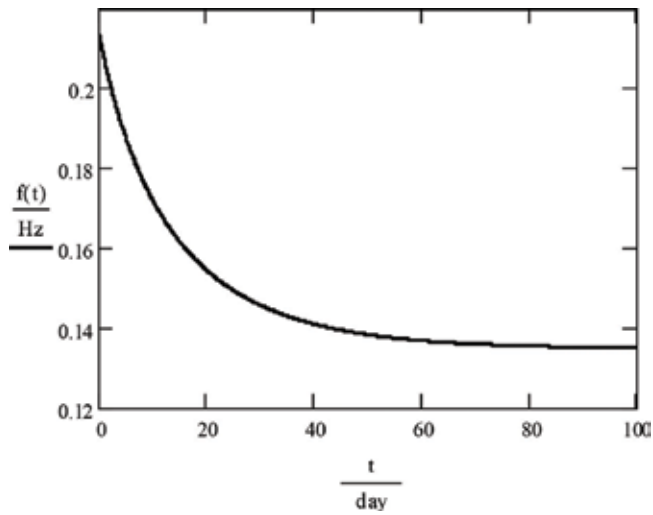


Figure 12. Frequency variation on structure at 90 days.

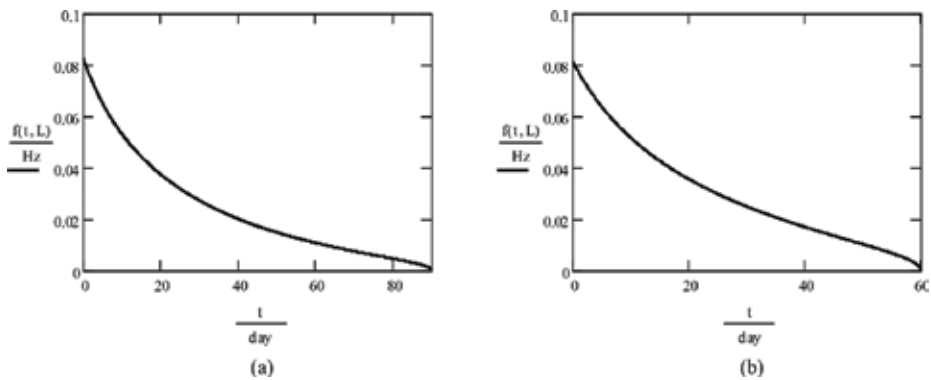


Figure 13. Structural frequency on height of losing stability. (a) $L = 50.6975$ m—Collapse on the 90th day, (b) $L = 57.0000$ m—Collapse on the 60th day.

5. Conclusions

Simulation 1

- In this work, a numerical simulation of a reinforced-prestressed concrete beam as a support for rotating machines was performed.
- For the vibration analysis, the viscoelasticity, which is an intrinsic material property, was introduced due to the slenderness of the beam, revealing a resonant regime not foreseen in the linear analysis (without viscoelasticity).

- The effect of geometric stiffness produced by the horizontal loading and the corresponding possibility of introducing resonant regimes in the structural support system were demonstrated by calculating their frequencies.
- It can be concluded, therefore, that due to the increase of the axial compressive force, resonance conditions can occur, as represented by the intersection of the curves in **Figure 7**. In the present study, resonance occurs when the axial compression force reaches 210 kN at 10 days. Other instants might also be considered.
- Since the force of post-tension decreases the stiffness of the beam, this can lead to the resonant regime if it has not been previously evaluated in the structural analysis.
- The technique studied in this chapter offers an efficient tool to provide the removal of the support structure of that unwanted regime, avoiding the production of harmful effects on the equipment, fabricated products, and work environment of the operators.
- In further work, it is necessary to introduce normative criteria, perform experimental activity, and evaluate the influence of the prestressing bar stiffness on the structural response.

Simulation 2

- The modulus of elasticity calculated by Eq. (6) on the ninth day was 16027.64 MPa, which represents a decrease of 49.81% in relation to the initial value of 31931.05 MPa.
- The frequency of the structure calculated at the initial moment was 0.215715 Hz, and on the 90th day, of 0.135021 Hz, representing a reduction of 37.41%.
- The simulated structure finds its limit of stability when reaching 50.6975 m, collapsing at 90 days. If the viscoelastic effect were not considered, the height limit would be 71.29 m, 20.47% above the first one. The obtained result was taken for an exactitude of five decimal significant algorithms ($f = 0.00000$).
- The previous aspect is relevant because if the height were considered between the limit established without the viscoelasticity and that defined with it, the structure would collapse before the end of 90 days in service. For a height of 57 m, for example, the collapse would occur 60 days after being loaded.
- Others rheological models for representing viscoelastic behavior can be tested in order to evaluate the frequency of a column of reinforced concrete as well as criteria from regulatory codes.
- The critical load of buckling can be obtained by using the same process present in this work and comparing it to other tools for calculations as, for example, finite element method (FEM).

Acknowledgements

This work was supported by the National Council for Scientific and Technological Development (CNPq) from Brazil by process 443,044/2014-7 in Call MCTI/CNPQ/Universal 14/2014.

Author details

Alexandre de M. Wahrhaftig^{1*}, Reyolando M. L. R. F. Brasil² and
Lázaro S. M. S. C. Nascimento¹

*Address all correspondence to: alixa@ufba.br

1 Department of Construction and Structures Rua Aristides Novis, Polytechnic School, Federal University of Bahia (UFBA), Salvador, BA, Brazil

2 Department of Structural and Geotechnical Engineering, University of São Paulo (USP) Polytechnic School, Cidade Universitária, São Paulo, SP, Brazil

References

- [1] Wahrhaftig AM, Brasil RMLRF, Balthazar JM. The first frequency of cantilever bars with geometric effect: a mathematical and experimental evaluation. *Journal of the Brazilian Society of Mechanical Sciences and Engineering*. 2013;**35**:457-467. DOI:10.1007/s40430-013-0043-9
- [2] Wahrhaftig AM, Brasil RMLRF. Representative experimental and computational analysis of the initial resonant frequency of largely deformed cantilevered beams. *International Journal of Solids and Structures*. 2016;**102-103**:44-55. DOI: 10.1016/j.ijsolstr.2016.10.018
- [3] Wahrhaftig AM, Brasil RMLRF. Vibration analysis of mobile phone mast system by Rayleigh method. *Applied Mathematical Modelling*. 2016;**42**:330-345. DOI: 10.1016/j.apm.2016.10.020
- [4] Wahrhaftig AM, Brasil RMLRF. Initial undamped resonant frequency of slender structures considering nonlinear geometric effects: The case of a 60.8 m-high mobile phone mast. *Journal of the Brazilian Society of Mechanical Sciences and Engineering*. 2016;**39**(3): 725-735. DOI: 10.1007/s40430-016-0547-1
- [5] Clough RW, Penzien J. *Dynamic of Structures*. 2nd ed. Taiwan: McGraw Hill International Editions; 1993
- [6] Findley WN, Lai JS, Onaran K. *Creep and Relaxation of Nonlinear Viscoelastic Materials, Whit an Introduction to Linear Viscoelasticity*. New York: Dover Publications, Inc.; 1989
- [7] Leohard F, Mong E. *Constructions of Concrete – Basic Principles for Dimensioning of Concrete Structures*. 1st. ed. Vol. 1, Rio de Janeiro: Livraria Interciência;1977
- [8] Mehta PK, Monteiro PJM. *Concrete: Structure, Proprieties and Materials*. São Paulo: PINI; 1994

- [9] Armijo NC, Paola MD, Pinnola FP. Fractional characteristic times and dissipated energy in fractional linear viscoelasticity. *Communications in Nonlinear Science and Numerical Simulation*. 2016;**37**:14-30. DOI:10.1016/j.cnsns.2016.01.003
- [10] Wahrhaftig AM, César SF, Brasil RMLRF. Creep in the fundamental frequency and stability of a slender wooden column of composite section. *Revista Árvore*. 2016;**40**(6):1129-1140. DOI:10.1590/0100-67622016000600018
- [11] Oliveira AW. Uso do concreto protendido em pontes: Estudo de Caso da Nova Ponte Sobre o Rio São Francisco na BR-101. V Encontro Nacional de Estudantes de Engenharia Civil, Porto Seguro; 2016. (In Portuguese)
- [12] Calduro EL. Em favor da leveza. *Revista Técnica, PINI*, n. 26; 1997. (In Portuguese)
- [13] Jost DT. Análise de peças fletidas com protensão não aderente pelo método dos elementos finitos. 152 p. Dissertação (Mestrado em Engenharia Civil) – Universidade Federal do Rio Grande do Sul. Porto Alegre, 2006. (In Portuguese)
- [14] Rayleigh. *Theory of Sound* (two volumes). New York: Dover Publications, re-issued;1877
- [15] Temple G, Bickley WG. *Rayleigh's Principle and its Applications to Engineering*. London: Oxford University Press, Humphrey Milford;1933
- [16] Leissa AW. The historical bases of the Rayleigh and Ritz methods. *Journal of Sound and Vibration*. Nov. 2005;**287**(4–5):961-978
- [17] Levy R, Spillers WR. *Analysis of Geometrically Nonlinear Structures*. New York: Chapman & Hall; 1995
- [18] National Regulatory Standard (Norma Brasileira), NBR – 6118–Design of Structural Concrete–Procedure. Rio de Janeiro: Brazilian National Standards Association (Associação Brasileira de Normas Técnicas, ABNT); 2014
- [19] Wahrhaftig AM. Analysis of the first modal shape: Using two case studies. *International Journal of Computational Methods*. 2018;**15**(3):1840019 (14 pages). DOI: 10.1142/S0219876218400194
- [20] Wahrhaftig AM, Brasil RMLRF. Essay on Creep in Vibration Columns (Ensaio sobre a fluência na vibração de colunas), Paper ID 002, Congress on Numerical Methods in Engineering CMN2015, Lisboa; 2015
- [21] Wahrhaftig AM. Rayleigh with Viscoelasticity Applied to a Highly Slender 40-m-high Concrete Mast, Congress on Numerical Methods in Engineering CMN2017, July 3–5, Valencial; 2017

Collapse Load for Thin-Walled Rectangular Tubes

Kenichi Masuda and Dai-Heng Chen

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.71226>

Abstract

In this chapter, thin-walled rectangular tubes under pure bending are considered, by performing a series of FEM numerical studies. In the simulation, a homogeneous and isotropic elastic perfectly plastic material was employed for the tube material. A commonly used method for predicting the collapse load of rectangular tubes subjected to pure bending was proposed by Kecman. Kecman's method focuses on a slenderness of the flange. When buckling occurs in the flange, this method uses a collapse load corresponding to the post buckling strength of the flange. When buckling does not occur at the flange, this method used a relation of the flange slenderness to the cross-sectional fully plastic yielding. This method for predicting the collapse loads is effective when the aspect ratio of web to flange is not large. However, for large aspect ratios, there is a large discrepancy between the values of maximum moment corresponding to the collapse loads obtained from this method and the FEM numerical results due to an effect of web slenderness. A new method is proposed to predict the maximum moment considering the effect of web slenderness. The validity of the collapse load estimation is checked by the results of FEM numerical simulation.

Keywords: thin-walled tube, bending, buckling, collapse load, fem

1. Introduction

The aims of this chapter are as follows:

- To understand the validity of existing estimation methods [1] by using the results of numerical simulations.
- To point out a case in which the estimation method is not applicable by using the results of numerical simulations.
- To understand a factor of the discrepancy by using the results of numerical simulations.

- To propose a new estimation method by considering the factor and using mathematical approach.
- To understand the validity of the new estimation method by comparing with the results of numerical simulations.

We selected “Collapse load for thin-walled rectangular tubes under bending” as the subject of these topics. The research content of this chapter is based on our recent paper [2], and this chapter shows the results of numerical simulations in detail.

2. To understand the validity of existing estimation methods by using the results of numerical simulations

2.1. Numerical simulation method

Figure 1(a) shows the simulated rectangular tubes, in which one end of the rectangular tube was fixed to a rigid wall, and pure bending was applied from the other end by modeling a lid rotating about the z axis under rotary control θ . Bending moment M can be derived from the rigid wall as reaction moment. **Figure 1(b)** shows a deformation shape and bending angle θ . Until buckling occurs, axial strain ε_x can be defined by

$$\varepsilon_x = \frac{\theta}{L}y \quad (1)$$

Figure 1(c) shows the axial strain distribution ε_x on cross-section for a square tube with $t = 0.4$ mm, $a = b = 50$ mm at $\theta/L = 0.01$ m⁻¹. As shown in **Figure 1(c)**, the axial strain distribution ε_x of FEM is in good agreement with the value obtained from Eq. (1). The effects of various geometric parameters were investigated under bending collapse. These parameters were tube thickness t , width of the flange a , and width of the web b . In order to prevent torsional behavior, a rigid lid was adopted as suggested by Guarracino [3]. In particular, the lid thickness t_f was set to five times of t . In the simulation, a homogeneous and isotropic elastic perfectly plastic material was employed for the tube material. As a yield condition, von Mises yield conditions were adopted. In this chapter, the material mechanical properties are set as follows. Young’s modulus E is set as 72.4 GPa, the yield stress σ_s is set as 72.4 MPa, and Poisson’s ratio ν is set as 0.3.

In this chapter, in order to formulate the geometric nonlinear behavior and solve the nonlinear equation, the updated Lagrange method, algorithm based on the Newton–Raphson method, and return-mapping method were used. The rectangular tubes were meshed using four-node quadrilateral thickness shell elements (Element type 75) with five integration points across the thickness. A convergence test on element size was conducted, and the adopted divide method was that the wall width divided into at least 20 sublengths, and the wall length divided as the elements become almost square.

In order to neglect the influence of the boundary conditions, the ratio of the length and width L/a , L/b was set to $L/a > 6$, $L/b > 6$. It means that the length of tubes was assumed to be large enough.

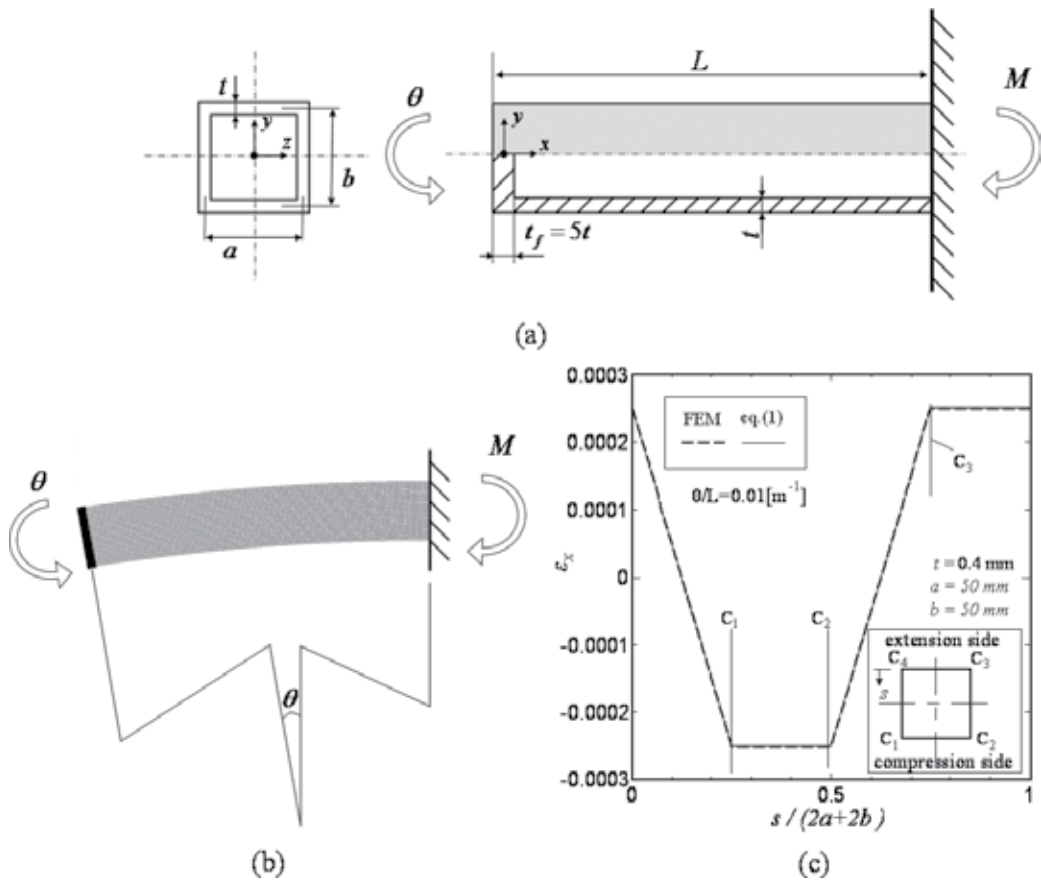


Figure 1. Numerical simulation model: (a) rectangular tube to which a pure bending moment is applied; (b) deformed shape; and (c) axial strain distribution on cross-section at $\theta/L = 0.01 \text{ m}^{-1}$.

2.2. Kecman’s method for predicting the maximum bending moment of rectangular tubes

Kecman focused on slenderness corresponding to buckling stress of the compression flange and proposed a formula to predict the collapse load or the maximum moment M_{\max} . Depending on the value of buckling stress $\sigma_{\text{buc-a}}$ of the compression flange

$$\sigma_{\text{buc-a}} = \frac{k_a \pi^2 E}{12(1-\nu^2)} \left(\frac{t}{a}\right)^2 \quad (2)$$

three cases are distinguished, as shown in **Figure 2**. In Eq. (2), k_a is the buckling coefficient, which Kecman assumed to be

$$k_a = 5.23 + 0.16 \frac{a}{b} \quad (3)$$

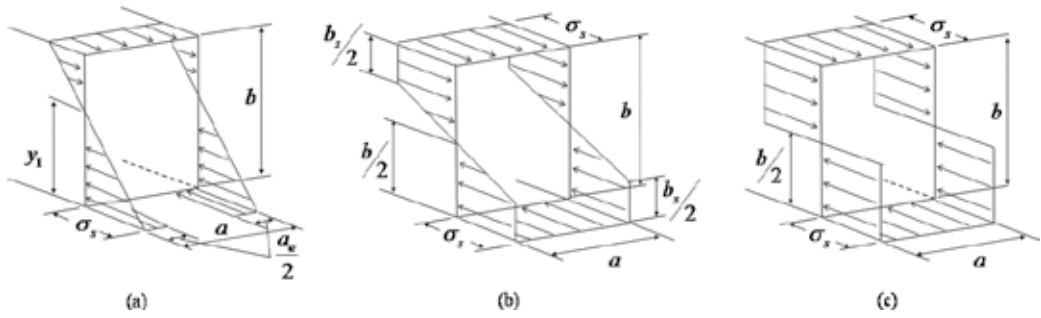


Figure 2. Schematic representation of axial stress distribution is used in the Kecman’s method: (a) case 1: $\sigma_{buc-a} < \sigma_s$; (b) case 2: $\sigma_s < \sigma_{buc-a} < 2\sigma_s$; and (c) case 3: $\sigma_{buc-a} > 2\sigma_s$.

The maximum moment M_{max} for the rectangular tube is given by

$$M_{max} = \sigma_s t b^2 \frac{2a + b + a_e \left(3 \frac{a}{b} + 2\right)}{3(a + b)} \tag{4}$$

For Case 1

$$M_{max} = M_{el} + (M_{pl} - M_{el}) \frac{\sigma_{buc-a} - \sigma_s}{\sigma_s} \tag{5}$$

For Case 2, and

$$M_{max} = M_{pl} \tag{6}$$

For Case 3. In the above equations

$$a_e = a \left(0.7 \frac{\sigma_{buc-a}}{\sigma_s} + 0.3\right) \tag{7}$$

and M_{el} and M_{pl} are the maximum elastic moment

$$M_{el} = \sigma_s t b \left(a + \frac{b}{3}\right) \tag{8}$$

and the cross-sectional fully plastic bending moment

$$M_{pl} = \sigma_s t b \left(a + \frac{b}{2}\right) \tag{9}$$

respectively.

Figure 3 shows a flow chart of the Kecman’s method for predicting the maximum moment of tubes under pure bending.

2.3. The applicability of the Kecman’s method for square tubes

Figure 4 shows that the bending moment M and the axial stress σ_x on cross-section for a square tube with $t = 0.4$ mm, $a = b = 50$ mm are subjected to pure bending ($\sigma_{buc-a} = 0.31\sigma_s < \sigma_s$). In order to better understand the bending collapse, Eq. (4) for Case 1 and the elastic buckling stress

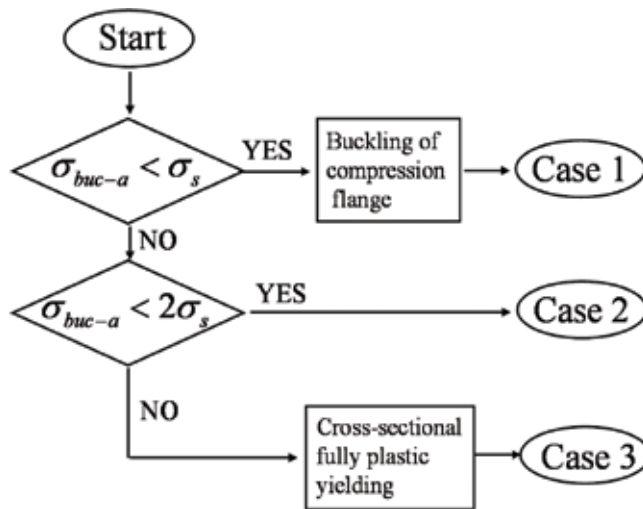


Figure 3. Flow chart of the Kecman's method for predicting the maximum moment of tubes under pure bending.

σ_{buc-a} given by Eq. (2) are also shown as a comparison. As shown in Figure 4(a), the maximum moment of FEM is in agreement with Eq. (4). The maximum value of σ_x at point P is in good agreement with Eq. (2). In addition, the axial compression stress σ_x at point Q at the quarter-web width keeps increasing after buckling occurs at point P in the middle of the compression flange, and the maximum value of σ_x at point Q occurs in the maximum moment. Moreover, as shown in Figure 4(b), although the axial compression stress in the middle of the compression flange decreases due to buckling at the flange, the axial compression stress increases at

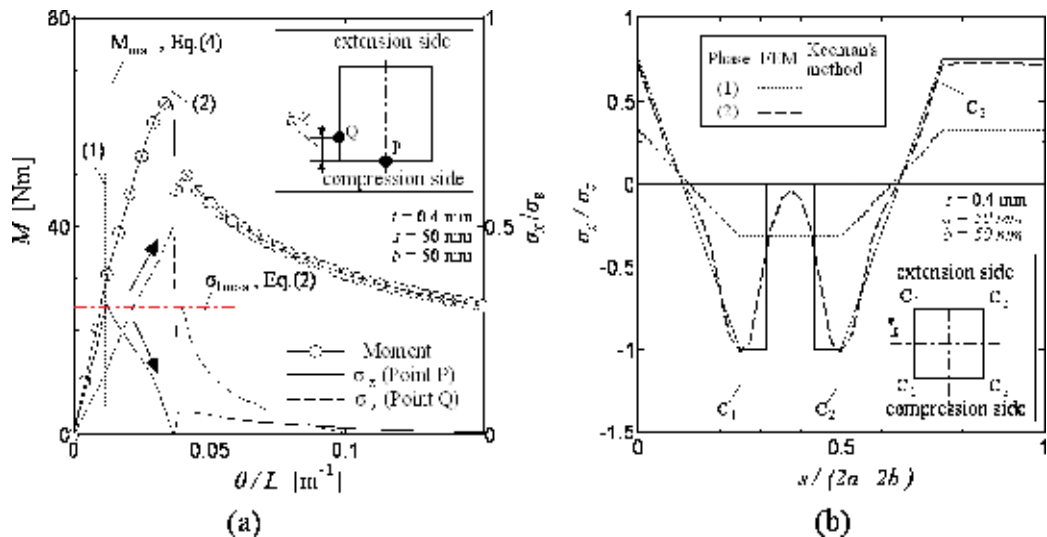


Figure 4. Moment M and axial stress σ_x on cross-section for a square tube with $t = 0.4$ mm, $a = b = 50$ mm are subjected to pure bending: (a) changes in moment M and axial stress σ_x at points P and Q on cross-section and (b) axial stress distribution on cross-section at phases (1) and (2) corresponding to $\theta/L = 0.012 m^{-1}$ and $\theta/L = 0.038 m^{-1}$, respectively, as denoted in (a).

both edges of the flange due to a corner constraint at the edges. Just after buckling, the stress increment at both edges is greater than the stress decrement in the middle of the compression flange, and thus the total force on the compression side increases and the moment increasing continuously. It is also noted that the stress on the web changes almost linearly; this suggests that buckling does not occur at the web. Therefore, the axial stress distribution when the maximum moment occurs is in good agreement with that obtained by the Kecman's method using the effective width of the compression flange, as shown by the solid line in the figure.

The above investigation confirms that for such tubes with $b/a = 1$ and $\sigma_{buc-a} < \sigma_s$, collapse is due to buckling at the compression flange, and the maximum moment can be predicted by the Kecman's method for Case 1.

Figure 5 shows the bending moment M and the axial stress σ_x on cross-section for a square tube with $t = 0.9$ mm, $a = b = 50$ mm ($\sigma_{buc-a} = 1.52\sigma_s > \sigma_s$). As shown in **Figure 5(a)**, the maximum moment is in good agreement with the value obtained from Eq. (5) for Case 2. The maximum value of σ_x at point P and Q occurs in the maximum moment and σ_x/σ_s at point P becomes 1. Moreover, as shown in **Figure 5(b)**, the absolute value of the axial stress at phase (2), for which the maximum moment occurs, is greater than the value at phase (1) for all cross-sectional positions. At phase (2), the stress at the flanges is equal to the yield stress σ_s and there also exist plastic yielding regions in the webs. In **Figure 5(b)**, Kecman's stress distribution when the maximum moment occurs is obtained by linear interpolation of two theoretical stress distributions corresponding to M_{el} and M_{pl} , respectively, and is shown by a solid line. It is seen from **Figure 5(b)** that the axial stress distribution when the maximum moment occurs obtained from numerical simulation is in good agreement with Kecman's stress distribution.

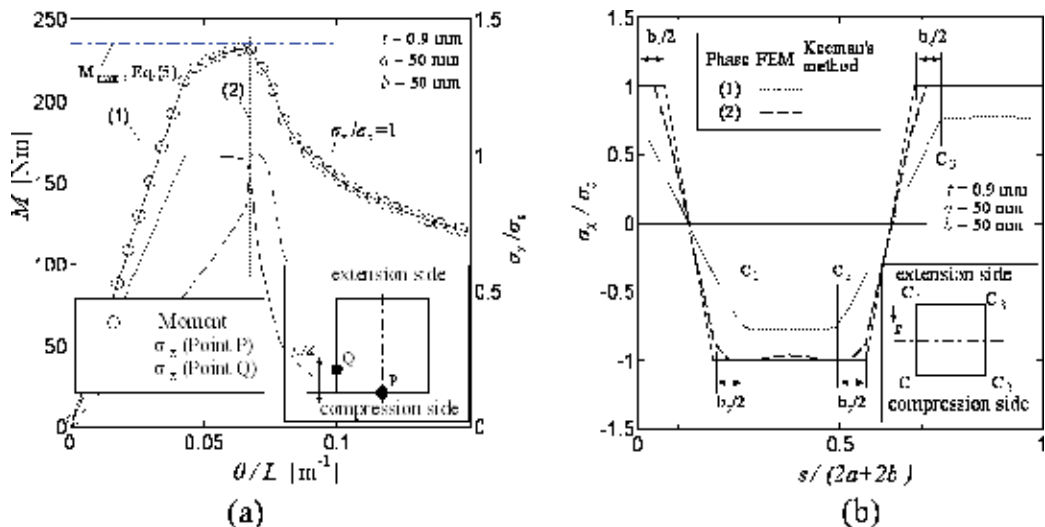


Figure 5. Moment M and axial stress σ_x on cross-section for a square tube with $t = 0.9$ mm, $a = b = 50$ mm are subjected to pure bending: (a) changes in moment M and axial stress σ_x at points P and Q on cross-section and (b) axial stress distribution on cross-section at phases (1) and (2) corresponding to $\theta/L = 0.025$ m⁻¹ and $\theta/L = 0.065$ m⁻¹, respectively, as denoted in (a).

The above investigation confirms that for such tubes with $b/a = 1$ and $\sigma_{buc-a} > \sigma_s$, the collapse is not due to buckling at the compression flange, but rather plastic yielding at the flange, and the maximum moment can be evaluated by Eq. (5) for Case 2.

3. To point out a case in which the estimation method is not applicable by using the results of numerical simulations

In order to investigate the accuracy of the Kecman's method for predicting the maximum moment M_{max} of tubes under bending, **Figure 6** shows the maximum bending moment of FEM numerical simulations for tubes with aspect ratios $b/a = 1, 2,$ and 3 . Eqs. (4), (5) and (6) are also shown as a comparison. As shown in the figure, the prediction of the Kecman's method is well in agreement with the results of FEM numerical simulations when the relative thickness t/a is not very small and the aspect ratio of web to flange b/a is not large, for example, when the tube relative thickness is about $t/a > 0.008$ for $b/a = 1$ and is about $t/a > 0.016$ for $b/a = 2$. However, for large aspect ratios, there is a large discrepancy between the values of maximum moment obtained from the Kecman's method and the FEM numerical results. This means that tubes with cross-section of a large aspect to which the Kecman's method does not apply are found

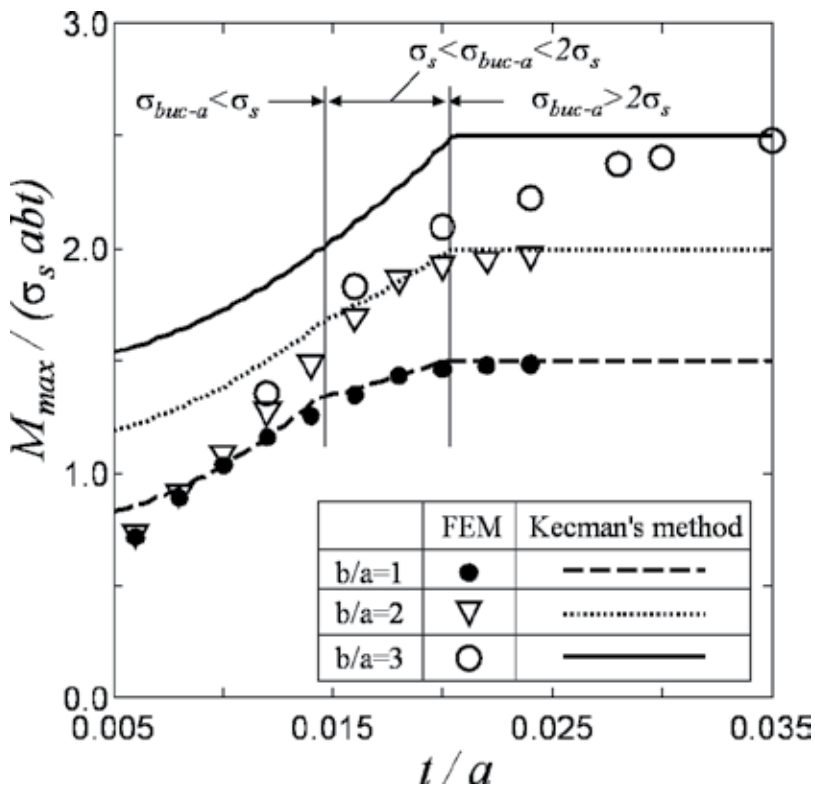


Figure 6. Comparison of the Kecman's method and the FEM numerical results.

to exist. Therefore, it is necessary to investigate the bending collapse mechanism of rectangular tubes in order to give an effective method for predicting the maximum moment of tubes.

4. To understand a factor of the discrepancy by using the results of numerical simulations

We investigate three tubes to which the Kecman's method is not applicable.

Figure 7 shows the bending moment M and the axial stress σ_x on cross-section for a rectangular tube with $t = 0.4$ mm, $a = 50$ mm, $b = 100$ mm ($\sigma_{\text{buc-a}} = 0.30\sigma_s$). As shown in **Figure 7(a)**, the maximum moment is less than the value obtained from Eq. (4) for Case 1. The maximum value of σ_x at point P is in good agreement with the elastic buckling stress $\sigma_{\text{buc-a}}$ given by Eq. (2), and the maximum value occurs before the maximum moment. Meanwhile, the axial compression stress σ_x at point Q at the quarter-web width decreases also before the moment becomes the maximum moment. Moreover, as shown in **Figure 7(b)**, the axial stress in the compression flange is concentrated at the edges when the maximum moment occurs, and the axial stress on the compression web does not change linearly. This suggests that compression buckling also arises at the web. Therefore, the axial stresses on the web at the maximum moment are less than that obtained by the Kecman's method, as indicated by the arrows in **Figure 7(b)**.

The above investigation reveals that, in cases when b/a are large and $\sigma_{\text{buc-a}} < \sigma_s$, the collapse is not only due to buckling at the compression flange but also due to buckling at the compression web. Therefore, the maximum moment cannot be predicted by the Kecman's method. Based on **Figure 7(b)**, the cross-sectional stress distribution under the maximum moment corresponding to this collapse mode can be schematically represented by **Figure 8(a)** and called Case 4 in this chapter.

Figure 9 shows the bending moment M and the axial stress σ_x on cross-section for a rectangular tube with $t = 0.5$ mm, $a = 20$ mm, and $b = 100$ mm ($\sigma_{\text{buc-a}} = 2.83\sigma_s$). As shown in **Figure 9(a)**, the maximum moment is less than that obtained from Eq. (6) for Case 3. The axial compression stress σ_x at point P in the middle of the compression flange increases until the moment becomes the maximum moment, and the value σ_x/σ_s becomes approximately 1. However, the axial compression stress σ_x at point Q at the quarter-web width decreases before the moment becomes the maximum moment. Also, as shown in **Figure 9(b)**, the axial stress distribution in the compression flange is almost constant, and the absolute value of σ_x/σ_s is approximately 1 when the maximum moment occurs. However, as compared with Case 3 shown in **Figure 2(c)**, it is found that although the buckling stress of the flange $\sigma_{\text{buc-a}}$ obtained from Eq. (2) is higher than twice the yielding stress, $\sigma_{\text{buc-a}} = 2.83\sigma_s > 2\sigma_s$, a plastic yielding region is not found in the web. Moreover, the axial stress in the web does not change linearly and decreases greatly in the compression portion of the web. This suggests that compression buckling arises at the web. Therefore, the axial stress distribution at the maximum moment differs greatly from that obtained by the Kecman's method, as indicated by the arrows in **Figure 9(b)**, because in the Kecman's method, the buckling of web is not taken into account.

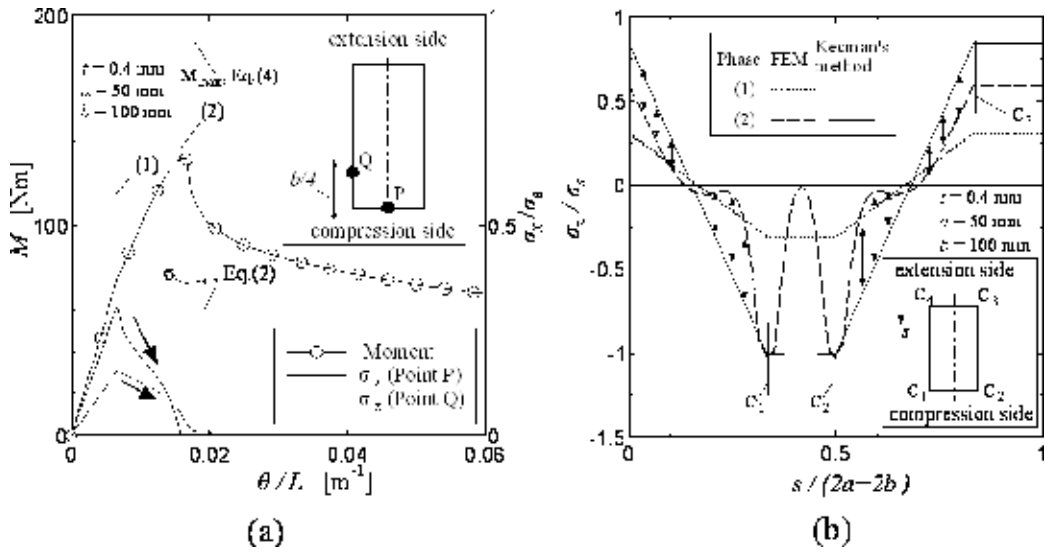


Figure 7. Moment M and axial stress σ_x on cross-section for a rectangular tube with $t=0.4$ mm, $a=50$ mm, $b=100$ mm are subjected to pure bending: (a) changes in moment M and axial stress σ_x at points P and Q on cross-section and (b) axial stress distribution on cross-section at phases (1) and (2) corresponding to $\theta/L=0.007$ m⁻¹ and $\theta/L=0.016$ m⁻¹, respectively, as denoted in (a).

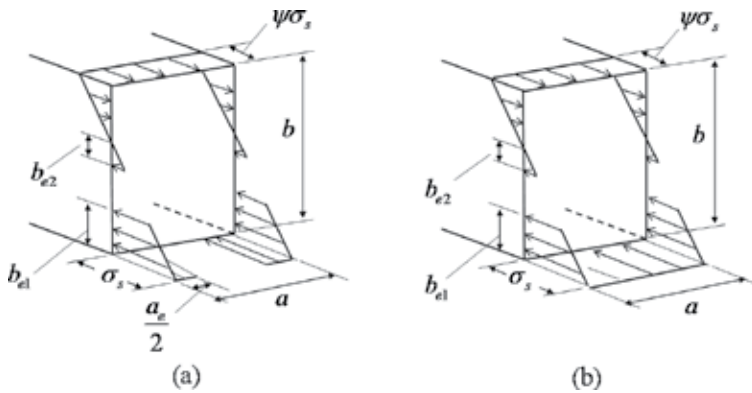


Figure 8. Schematic representation of axial stress distribution with considering the buckling at web when the maximum moment occurs: (a) case 4: $\sigma_{buc-a} < \sigma_s$ and (b) case 5: $\sigma_{buc-a} > \sigma_s$.

The above investigation reveals that, in such tubes with large aspect ratio b/a , even though $\sigma_{buc-a} > \sigma_s$, collapse is not due to plastic yielding at the flange, but rather buckling at the compression web in a state of plastic yielding at the compression flange. Therefore, the maximum moment cannot be predicted by the Kecman's method. Based on **Figure 9(b)**, the cross-sectional stress distribution under the maximum moment corresponding to this collapse mode can be schematically represented by **Figure 8(b)** and called Case 5 in this chapter.

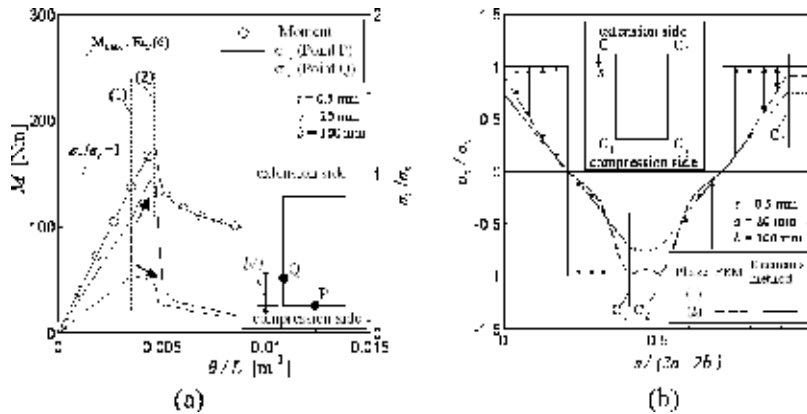


Figure 9. Moment M and axial stress σ_x on cross-section for a rectangular tube with $t = 0.5$ mm, $a = 20$ mm, $b = 100$ mm are subjected to pure bending: (a) changes in moment M and axial stress σ_x at points P and Q on cross-section and (b) axial stress distribution on cross-section at phases (1) and (2) corresponding to $\theta/L = 0.036$ m⁻¹ and $\theta/L = 0.048$ m⁻¹, respectively, as denoted in (a).

Figure 10 shows the bending moment M and the axial stress σ_x on cross-section for a rectangular tube with $t = 1.2$ mm, $a = 50$ mm, $b = 150$ mm ($\sigma_{buc-a} = 2.8\sigma_s$). As shown in **Figure 10(a)**, the maximum moment is less than that obtained from Eq. (6) for Case 3. The axial compression stress σ_x at point P in the middle of the compression flange increases up to the yielding stress σ_s before the maximum moment was reached and sets the value σ_x/σ_s equal approximately to 1 until the moment becomes the maximum moment. However, the axial compression stress σ_x at point Q at the quarter-web width increases until the moment becomes the maximum

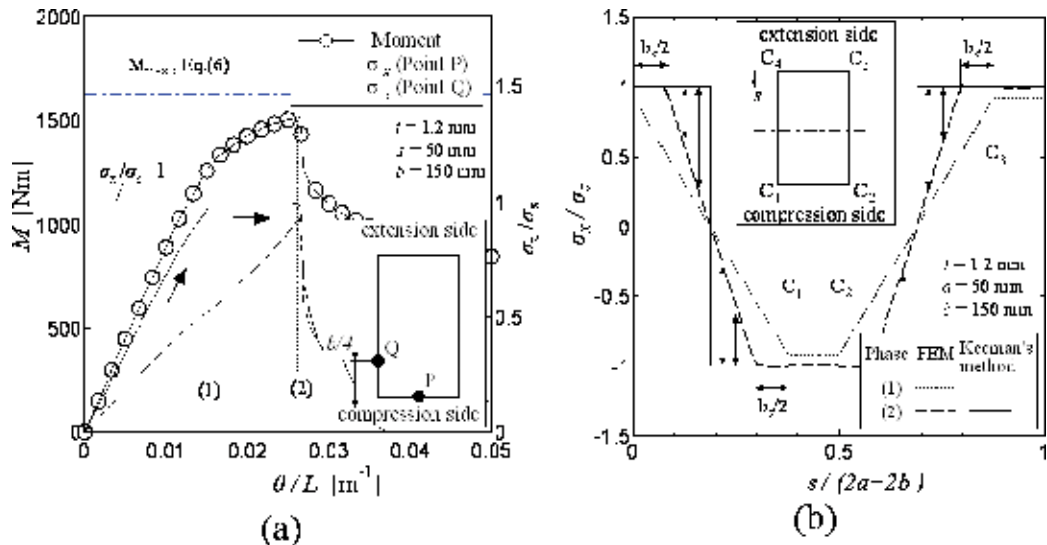


Figure 10. Moment M and axial stress σ_x on cross-section for a rectangular tube with $t = 1.2$ mm, $a = 50$ mm, $b = 150$ mm are subjected to pure bending: (a) changes in moment M and axial stress σ_x at points P and Q on cross-section and (b) axial stress distribution on cross-section at phases (1) and (2) corresponding to $\theta/L = 0.015$ m⁻¹ and $\theta/L = 0.026$ m⁻¹, respectively, as denoted in (a).

moment. Also, it is seen from **Figure 10(b)** that the axial stress distribution in the compression flange is almost constant, and the absolute value of σ_x/σ_s is approximately 1 when the maximum moment occurs. Moreover, it is also found from a comparison with **Figure 9(b)** that in the web, no buckling occurs, but plastic yielding regions can be observed. However, the plastic yielding is not generated to the entire web, although the buckling stress of the flange $\sigma_{\text{buc-a}}$ obtained from Eq. (2) is higher than twice the yielding stress, $\sigma_{\text{buc-a}} = 2.8\sigma_s > 2\sigma_s$. Thus, the stress distribution is different from the cross-sectional fully plastic yielding, as indicated by the arrows in the figure. This suggests that even if a compression buckling does not arise at the web, the web slenderness also affects the cross-sectional fully plastic yielding of the tube under bending. That is, the conditions of generating the cross-sectional fully plastic yielding are dependent not only on the flange slenderness but also on the web slenderness. In the Kecman's method, the conditions for the cross-sectional fully plastic yielding are determined by only the ratio of $\sigma_{\text{buc-a}}$ to σ_s .

The above investigation reveals that in such tubes with large aspect ratio b/a , even though $\sigma_{\text{buc-a}} > 2\sigma_s$, the cross-sectional stress distribution under the maximum moment corresponding to this collapse mode may differ from that of the cross-sectional fully plastic yielding. Therefore, the maximum moment for such tubes cannot be predicted by the Kecman's method.

5. To propose a new estimation method by considering the factor and using mathematical approach

5.1. Effect of the web slenderness on the buckling at web

Bending stress occurs in the web of tube. The problem of web buckling is expressed in **Figure 11**. In **Figure 11(a)**, plate ABCD is defined by the width b and thickness t . As a boundary condition, displacement in the out-of-plane direction (displacement in the z direction) is fixed at both longitudinal edges (BC and DA). The bending and compression are applied through displacement control. For the ultimate loading after buckling, the distribution of compressive stress σ_x along the width direction is characterized by two effective widths, b_{e1} and b_{e2} , as shown in **Figure 11(b)**. In the figure, compressive stress is denoted by a positive value.

Many studies have been reported on the ultimate loading of a plate after buckling under bending and compression. For example, the effective widths b_{e1} and b_{e2} for a plate under stress gradient shown in **Figure 11** are given in AS/NZS 4600 standard [4] and NAS [5] as follows:

$$\begin{cases} b_{e1} = \frac{b_e}{3-\psi} \\ b_{e2} = \begin{cases} b_e/2 & \text{when } \psi \leq -0.236 \\ b_e - b_{e1} & \text{when } \psi > -0.236 \end{cases} \end{cases} \quad (10)$$

In addition, $b_{e1} + b_{e2}$ shall not exceed the compression portion of the web. Here, ψ is ratio of f_2^* and f_1^* , f_2^* , f_1^* and f_2 are web stresses shown in **Figure 11(b)**.

$$\psi = \frac{f_2^*}{f_1^*} \quad (11)$$

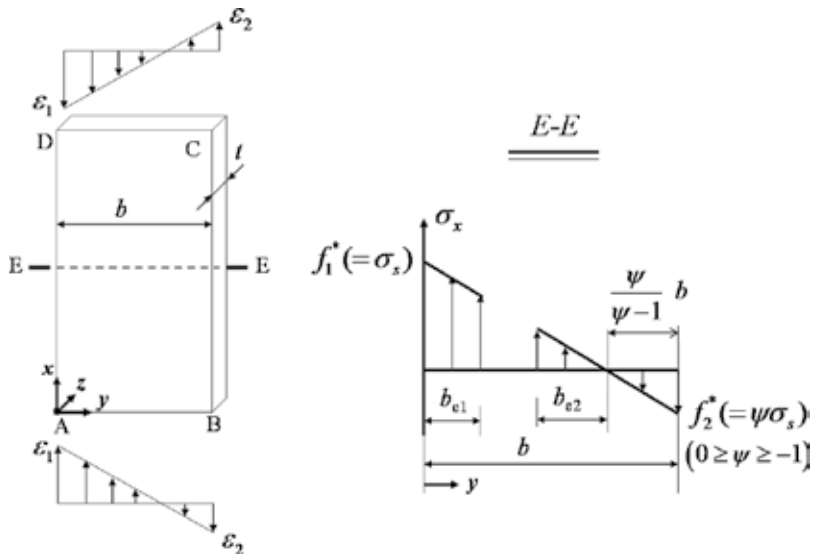


Figure 11. Plate subjected to compression and bending: (a) analyzed model and (b) axial compressive stress σ_x distribution on E-E cross-section in (a).

λ is defined by

$$\lambda = \sqrt{\frac{\sigma_s}{\sigma_{buc-b}}} \tag{12}$$

The elastic buckling stress of web σ_{buc-b} is calculated as follows:

$$\sigma_{buc-b} = \frac{k_b \pi^2 E}{12(1-\nu^2)} \left(\frac{t}{b}\right)^2 \tag{13}$$

where the buckling coefficient k_b is given by

$$k_b = 4 + 2(1-\psi)^3 + 2(1-\psi) \tag{14}$$

b_e is given by

$$b_e = \rho b \tag{15}$$

ρ is called the reduction factor and is given by

$$\rho = \frac{1}{\lambda} \tag{16}$$

which is proposed by von Karman et al. [6]. The following formula for ρ :

$$\rho = \frac{1}{\lambda} \left(1 - \frac{0.22}{\lambda}\right) \tag{17}$$

is also proposed by Winter [7] and is well used for design specifications. The reason of Eq. (15) modified to Eq. (17) in actual design is mainly due to the fact that the maximum load capacity

of a buckling plate is reduced greatly by imperfections when the buckling stress is close to the yield stress [8]. Therefore, Eq. (16) is desirable for the present model because the influence of imperfections is not taken into consideration here. Moreover, in order to consider continuity of the load capability of a web with $\lambda = 1$, for which elastic buckling does not occur because $\sigma_{\text{buc-b}} = \sigma_y$, we apply Eq. (16) to the present study.

Eq. (10) is applied to the webs investigated in **Figures 7(b)** and **9(b)** to determine the corresponding effective width; the stress distributions on the web based on the obtained effective width using Eq. (10) are shown in **Figures 12(a)** and **13(a)**. In **Figure 12(a)**, the stress distribution obtained using Eq. (10) is qualitatively corresponding with the redistribution of the compression stress after buckling obtained from the FEM numerical simulation. However, in **Figure 13(a)**, even though there is a fall of the compression stress in the compression portion of the web after buckling as shown by the FEM simulation, the stress distribution obtained from Eq. (10) looks like a straight line because the effective widths b_{e1} and b_{e2} determined by Eq. (10) satisfy the following equation:

$$b_{e1} + b_{e2} = \frac{b}{1-\psi} \tag{18}$$

which means $b_{e1} + b_{e2}$ is equal to the compression portion of the web.

In fact, when the effective width is determined using Eq. (10), there are many instances in which Eq. (18) is satisfied. **Figure 14** shows various possible values of buckling stress of web, for which Eq. (18) is satisfied, for various assumed stress ratios ψ by solid line, as evaluated in Eq. (10) with ρ defined in Eq. (16). In **Figure 14**, the dashed line shows the corresponding result if ρ was calculated using Eq. (17); it is also seen from the dashed line that even if Eq. (17) is used for ρ the instances in which Eq. (18) is satisfied still exist. For these instances, the redistribution of the compression stress after buckling cannot be expressed by the effective width obtained from Eq. (10); this means that there is a possibility of giving a too large load

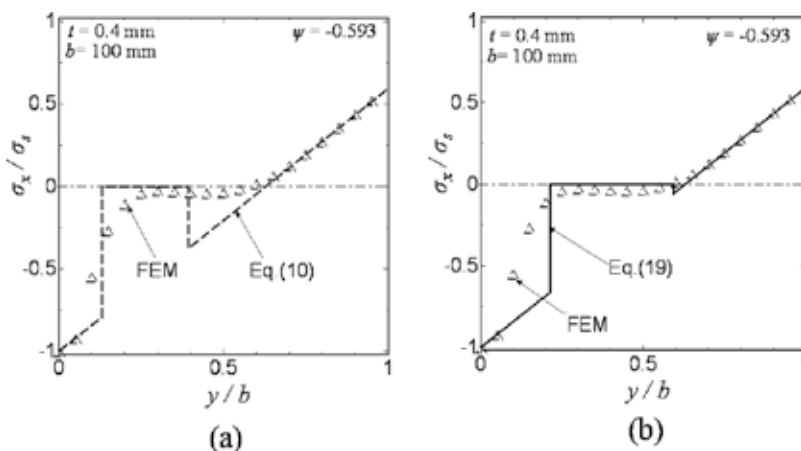


Figure 12. Stress distribution of web when the ultimate load is reached for the tube used in **Figure 7(b)**: (a) comparison with Eq. (10) and (b) comparison with Eq. (19).

capability of web from Eq. (10). Therefore, here as a comparison, we also use another solution given by Rusch and Lindner [9] which is given for the same plate shown in **Figure 11(a)** but with one of the two longitudinal edges BC being free. Although the free boundary condition at the longitudinal edge BC is different from the actual situation of web constituting the tube, the effect is assumed to be small because BC is under tension stress.

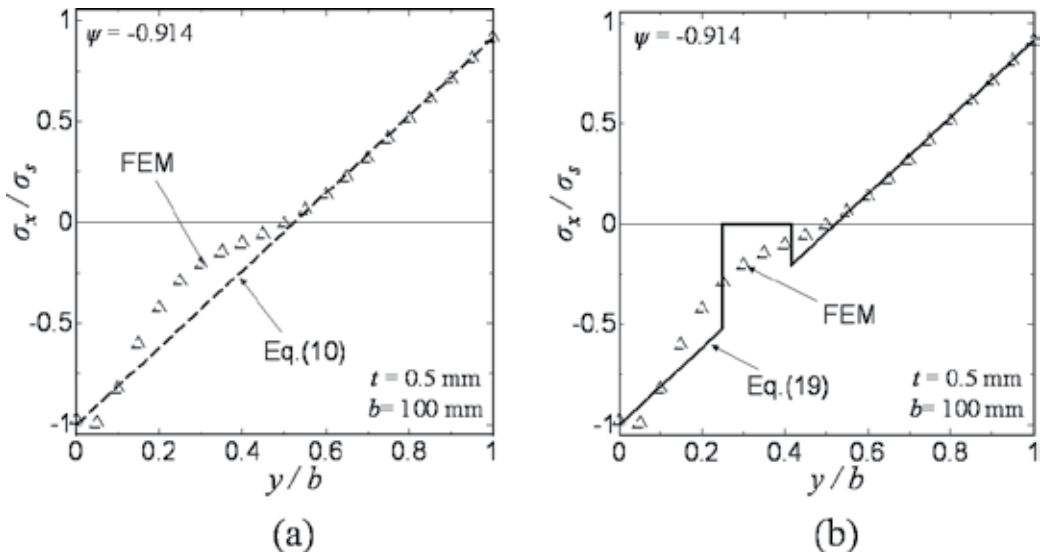


Figure 13. Stress distribution of web when the ultimate load is reached for the tube used in **Figure 9(b)**: (a) comparison with Eq. (10) and (b) comparison with Eq. (19).

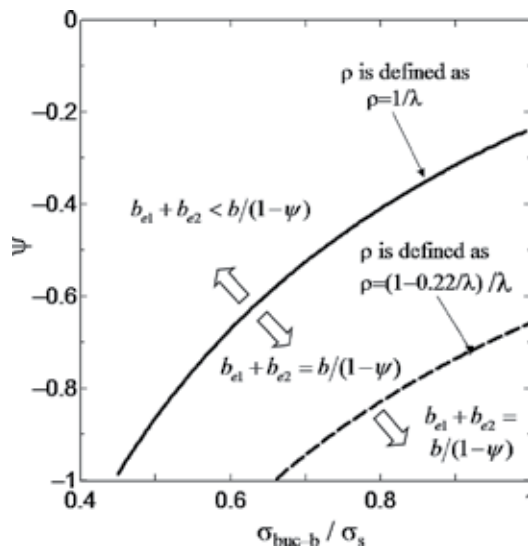


Figure 14. Various possible buckling stress σ_{buc-b} and stress ratio ψ with Eq. (18) satisfied.

In Ref. [9], the effective widths b_{e1} and b_{e2} are given by

$$\begin{cases} b_{e1} = b_e - b_{e2} \\ \frac{b_{e2}}{b} = \frac{0.226}{\lambda^2} \end{cases} \quad (19)$$

where

$$\frac{b_e}{b} = \frac{\rho}{1 - \psi} \quad (20)$$

Here, λ and ρ are calculated by Eqs. (12) and (16), respectively, the buckling stress $\sigma_{\text{buc-b}}$ is determined by Eq. (13) with k_b determined as follows:

$$k_b = 1.7 - 5\psi + 17.1\psi^2 \quad (21)$$

Figures 12(b) and **13(b)** show the comparisons of stress distributions on the web obtained from FEM and Eq. (19) for the tubes used in **Figures 7(b)** and **9(b)**, from which it is seen that the redistribution of stress after web buckling can be approximately expressed using Eq. (19). Comparing (a) and (b) in **Figure 12**, it is seen that for the stress distribution on the web in the tube used in **Figure 7(b)**, Eq. (19) is inferior in accuracy to Eq. (10). However, as shown in **Figure 13**, which shows the stress distributions on the web for the tube used in **Figure 9(b)**, although the fall of the compression stress in the compression portion of the web after buckling is not expressed by the solution obtained from Eq. (10), it is expressed by the solution from Eq. (19). In fact, it is seen from Eq. (20) that for the stress distribution on the web as obtained from Eq. (18), the length of $b_{e1} + b_{e2}$ is always smaller than the compression portion of the web.

5.2. Effect of the web slenderness on the cross-sectional fully plastic yielding

For tubes with large aspect ratio of web to flange, as an effect of web slenderness on the tube collapse, we considered the possible buckling of web and thus investigated the existence of Cases 4 and 5, as shown above. Hereafter, we consider the other effect of web slenderness on the cross-sectional fully plastic yielding.

As shown in **Figure 6**, for tubes with $b/a = 3$, there is a large discrepancy between Kecman's prediction and the FEM simulation. When the tubes are very thin (e.g., when $t/a < 0.02$ for $b/a = 3$) it is thought that the error generating is brought about because the web buckling was not taken into consideration in the Kecman's method. However, for the relatively thick tubes, the cause which produces the error is clearly different because buckling does not occur in such tubes. For example, for the tube with $b/a = 3$ and $t/a = 0.024$ shown in **Figure 10(b)**, even though the buckling stress of flange $\sigma_{\text{buc-a}}$ calculated by Eq. (2) is $\sigma_{\text{buc-a}}/\sigma_s = 2.8 > 2$, the maximum moment M_{max} as evaluated by FEM numerical simulation is $M_{\text{max}}/M_{pl} \approx 0.9$, which is not in agreement with Eq. (6) for the case of $\sigma_{\text{buc-a}}/\sigma_s = 2$ in the Kecman's method. Here, buckling does not occur in the web either because $\sigma_{\text{buc-b}}/\sigma_s = 1.4 > 1$. Also, it is seen from **Figure 10(b)** that the stress distribution on the cross-section is different from that shown in **Figure 2(c)** for Case 3 corresponding to the cross-sectional fully plastic yielding. This fact means that the condition for reaching the cross-sectional fully plastic yielding is also related to the web slenderness.

In order to consider the effect of the web slenderness on the tube collapse, the condition of $\sigma_{buc-a} > 2\sigma_s$ for Case 3 or for $M_{max} = M_{pl}$ in the Kecman’s method is replaced in the present study by the following condition:

$$\begin{cases} \sigma_{buc-a} \geq 2\sigma_s \\ \sigma_{buc-b} \geq 2\sigma_s \end{cases} \tag{22}$$

Here, σ_{buc-b} is determined assumed $\psi = -1$. When Eq. (22) is not satisfied, that is, when

$$\begin{cases} \sigma_s < \sigma_{buc-a} < 2\sigma_s \\ \sigma_{buc-b} \geq 2\sigma_s \end{cases} \tag{23}$$

or

$$\begin{cases} \sigma_{buc-a} \geq 2\sigma_s \\ \sigma_s < \sigma_{buc-b} < 2\sigma_s \end{cases} \tag{24}$$

or

$$\begin{cases} \sigma_s < \sigma_{buc-a} < 2\sigma_s \\ \sigma_s < \sigma_{buc-b} < 2\sigma_s \end{cases} \tag{25}$$

the stress on cross-section is expressed by Case 2 shown in **Figure 2(b)**. This fact can be confirmed from **Figure 10(b)** for which Eq. (24) is satisfied.

It is seen from the cross-sectional stress distribution shown in **Figure 10(b)** that the maximum moment in this case is dependent on the plastic yielding region in the web. Denoting the length of this plastic yielding region by b_s (see **Figure 2(b)**), the maximum moment can be evaluated through the value of b_s as follows:

for Case 2,

$$M_{max} = \sigma_s t \left[\frac{1}{6}(2b^2 + 2bb_s - b_s^2) + ab \right] \tag{26}$$

Substituting Eqs. (8), (9), and (26) into Eq. (5), b_s is obtained as

$$\frac{b_s}{b} = \begin{cases} 1 - \sqrt{2 - \left(\frac{t}{t_{ea}}\right)^2} & (\text{for } t_{ea} < t < \sqrt{2} t_{ea}) \\ 1 & (\text{for } t \geq \sqrt{2} t_{ea}) \end{cases} \tag{27}$$

where t_{ea} is the flange thickness for which the elastic buckling stress σ_{buc-a} obtained from Eq. (2) is equal to the yielding stress σ_s and is given by

$$t_{ea} = a \sqrt{\frac{12(1-\nu^2)}{k_g \pi^2}} \cdot \sqrt{\frac{\sigma_s}{E}} \tag{28}$$

Eq. (26) means that the b_s is determined by the flange slenderness only when Eq. (23) is satisfied. Therefore, when Eq. (24) is satisfied, we also suppose that the b_s can be determined by the web slenderness only as follows:

$$\frac{b_s}{b} = \begin{cases} 1 - \sqrt{2 - \left(\frac{t}{t_{eb}}\right)^2} & (\text{for } t_{eb} < t < \sqrt{2} t_{eb}) \\ 1 & (\text{for } t \geq \sqrt{2} t_{eb}) \end{cases} \quad (29)$$

where t_{eb} is the web thickness for which the elastic buckling stress σ_{buc-b} is equal to the yielding stress σ_s and is given by

$$t_{eb} = b \sqrt{\frac{12(1-\nu^2)}{k_a \pi^2}} \sqrt{\frac{\sigma_s}{E}} \quad (30)$$

Furthermore, we assume that this technique can also be used to evaluate the maximum moment in the case when Eq. (25) is satisfied. That is, M_{max} is determined from the smaller one from both values of b_s given in Eq. (27) and in Eq. (29). Validity of this assumption can be understood from **Figure 16(b)** shown later, in which for the tube with $t/a = 0.016$, Eq. (25) is satisfied: $\sigma_{buc-a} = 1.23\sigma_s$ and $\sigma_{buc-b} = 1.39\sigma_s$. Therefore, when using Eq. (26) to determine M_{max} for Case 2, the value of b_s is calculated using Eq. (27) if

$$t_{ea} > t_{eb} \quad (31)$$

and is calculated using Eq. (29) if

$$t_{ea} < t_{eb} \quad (32)$$

Using t_{ea} and t_{eb} , the condition of Eq. (22) can be rewritten as.

$$t \geq \sqrt{2} t_{ea} \text{ and } t \geq \sqrt{2} t_{eb} \quad (33)$$

5.3. Estimation of collapse load for thin-walled rectangular tubes under bending

Figure 15 shows a flow chart of a new method proposed in the present study for predicting the maximum moment of tubes under pure bending. This method includes both the possible buckling at web and the effect of web slenderness on the cross-sectional fully plastic yielding. In the flow chart, $\sigma_{buc-b,1}$ and $\sigma_{buc-b,2}$ are the buckling stress of web assuming the stress ratio ψ to be

$$\psi = -\frac{b - y_1}{y_1} = -\frac{a_e + b}{a + b} \quad (34)$$

and $\psi = -1$, respectively. Moreover, it is notable that in calculating the maximum bending moment for Cases 4 and 5 the stress ratio ψ is also unknown, which shall be determined from the conditions of pure bending through trial and error. Using the determined value of ψ , the maximum moment for Cases 4 and 5 is calculated as follows:

for Case 4:

$$\frac{M_{max}}{\sigma_s t} = \psi [b^2 + d_1^2 - d_2^2] + \frac{2(1-\psi)}{3b} [b^3 + d_1^3 - d_2^3] + a_e b \quad (35)$$

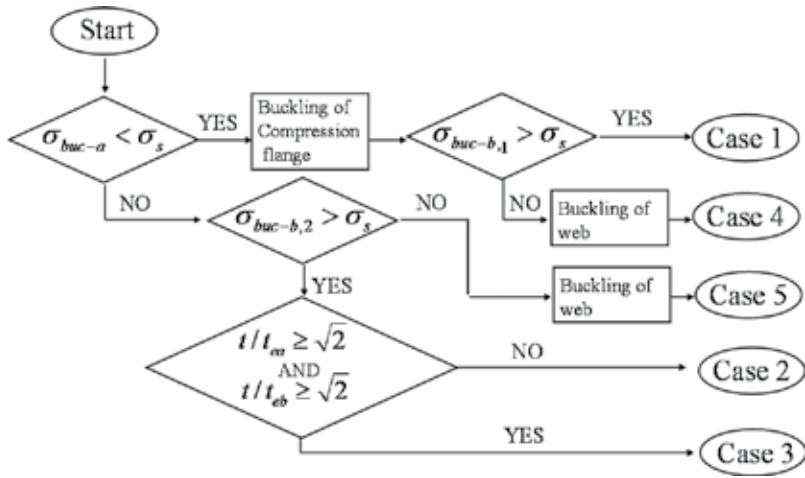


Figure 15. Flow chart of a new method proposed in the present study for predicting the maximum moment of tubes under pure bending.

for Case 5:

$$\frac{M_{\max}}{\sigma_s t} = \psi [b^2 + d_1^2 - d_2^2] + \frac{2(1-\psi)}{3b} [b^3 + d_1^3 - d_2^3] + ab \tag{36}$$

In Eqs. (35) and (36).

$$d_1 = b_{e2} + \frac{\psi}{\psi-1} b, d_2 = b - b_{e1} \tag{37}$$

6. To understand the validity of the new estimation method by comparing with the results of numerical simulations

In **Figure 16**, the maximum moment predicted by the present method is compared with that obtained from the FEM numerical simulation with $a/b = 1, 2,$ and 3 for **Figures 16(a), 16(b),** and **16(c)**, respectively. As a result of the prediction method proposed in this chapter, “method 1” uses Eq. (10) and “method 2” uses Eq. (19) for calculating the effective width.

The case number of the collapse corresponding to each thickness is also shown in the figures. In Case 2, there are two possible subcases: (1) $t_{ea} > t_{eb}$ as shown in **Figure 16(a)** and **(b)** and (2) $t_{ea} < t_{eb}$ as shown in **Figure 16(c)**; the maximum moment is determined by Eq. (27) for the former and by Eq. (29) for the latter. As shown in these figures, Eqs. (27) and (29) give good prediction to the corresponding subcase, respectively.

For Cases 4 and 5, although each result obtained from methods 1 and 2 is approximately in agreement with the FEM results of numerical simulations, it is found that there is a gap in the results between methods 1 and 2. When the buckling stress of the web σ_{buc-b} is close to

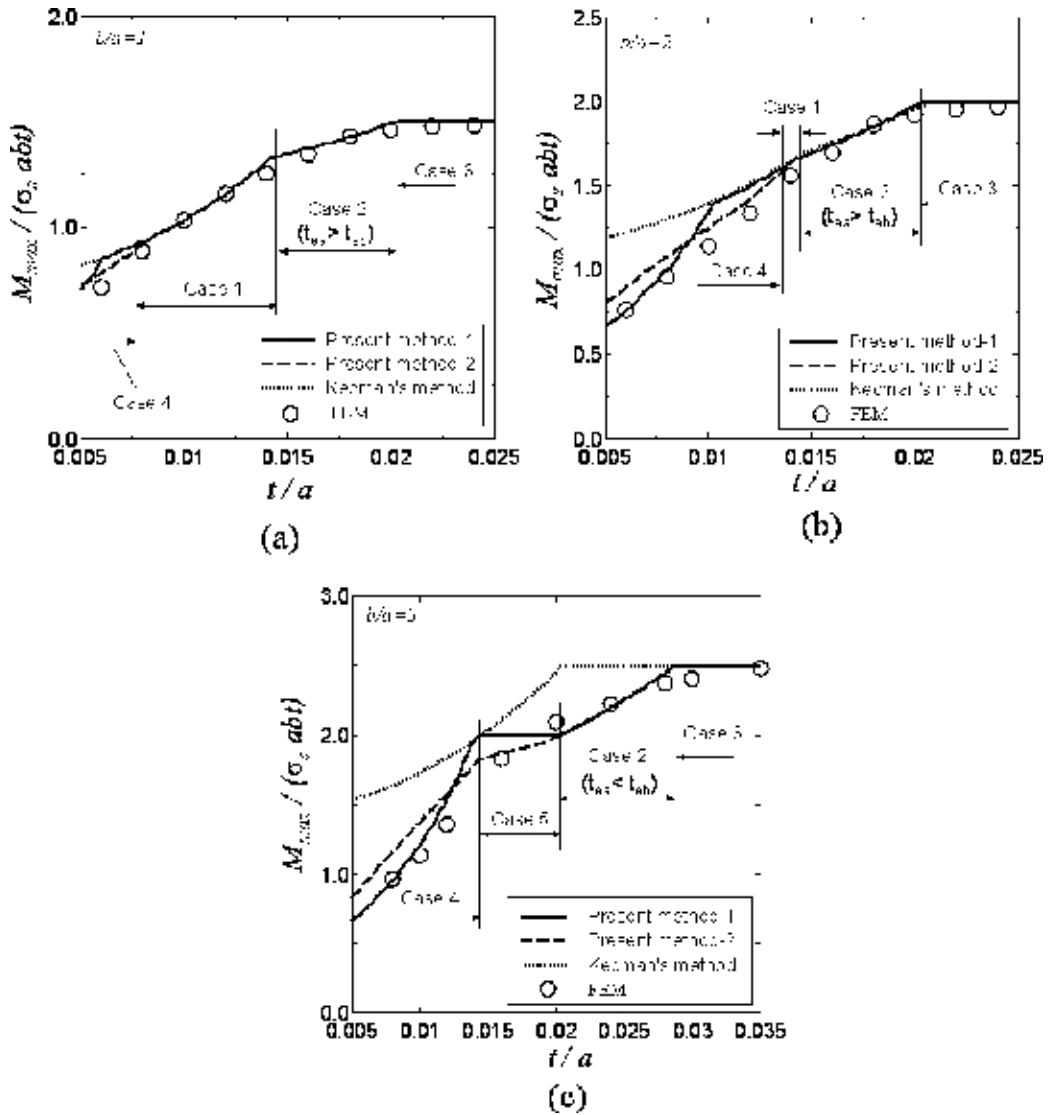


Figure 16. Prediction of the maximum bending moment M_{max} for rectangular tubes: (a) $a/b = 1$; (b) $a/b = 2$; and (c) $a/b = 3$.

the yielding stress σ_y , the method 1 gives a too large prediction as compared with the FEM results, reflecting the fact that $b_{e1} + b_{e2}$ given by Eq. (10) may be equal to the compression portion of the web as shown in Figure 14. However, for small t/a , when σ_{buc-b} is very much less than the yielding stress, method 1 is more accurate compared with method 2. Combining the advantages of these two methods, it is seen from Figure 16(a), (b), and (c) that the smaller one from both solutions obtained from method 1 and obtained from method 2 is in good agreement with the FEM results for all of Cases 4 and 5.

7. Conclusion

In this chapter, bending collapse of rectangular tubes was investigated using the FEM numerical simulation. The Kecman's method in which the post buckling strength of the flange and the effect of the flange slenderness on the cross-sectional fully plastic yielding are taken into account is effective when the aspect ratio of web to flange is not large. However, in order to predict accurately the maximum moments of rectangular tubes with large aspect ratio of web to flange, the slenderness of web has to be taken into account. Our new method in which the post buckling strength of the web under stress gradients and the effect of the web slenderness on the cross-sectional fully plastic yielding are taken into account are proposed, and the predicted maximum moment agrees with the results of FEM numerical simulations.

Author details

Kenichi Masuda^{1*} and Dai-Heng Chen²

*Address all correspondence to: masuda@eng.u-toyama.ac.jp

¹ Faculty of Engineering, University of Toyama, Toyama, Japan

² National Center for International Research on Structural Health Management of Critical Components, Jiangsu University, Zhenjiang, China

References

- [1] Kecman D. Bending collapse of rectangular and square section tubes. *International Journal of Mechanical Sciences*. 1983;**25**(9-10):623-636
- [2] Chen HD, Masuda K. Estimation of collapse load for thin-walled rectangular tubes under bending. *ASME Journal of Applied Mechanics*. 2016;**83**:031012-1-031012-8. DOI: 10.1115/1.4032159
- [3] Guarracino F. On the analysis of cylindrical tubes under flexure: Theoretical formulations, experimental data and finite element analyses. *Thin-Walled Structures*. 2003;**41**:127-147
- [4] AS/NZS 4600, editor. Cold-formed steel structures. Australian/New Zealand standard: Standards Australia; 2005
- [5] NAS, editor. AISI Standard. 2004 supplement to the North American specification for the design of cold-formed steel structural members 2001 edition. American Iron and Steel Institute; 2004

- [6] Karman VT, Sechler EE, Donnell HL. Strength of thin plates in compression. The American Society of Mechanical Engineers. 1932;**54**:53-57
- [7] Winter G. Strength of thin steel compression flanges. Transactions of the American Society of Civil Engineers. 1947;**112**:527-554
- [8] Rhodes J. Buckling of thin plates and members-and early work on rectangular tubes. Thin-Walled Structures. 2002;**40**(2):87-108
- [9] Rusch A, Lindner J. Application of level I interaction formulae to class 4 sections. Thin-Walled Structures. 2004;**42**(2):279-293

Edited by Srinivas P. Rao

Computational science is one of the rapidly growing multidisciplinary fields. The high-performance computing capabilities are utilized to solve and understand complex problems. This book offers a detailed exposition of the numerical methods that are used in engineering and science. The chapters are arranged in such a way that the readers will be able to select the topics appropriate to their interest and need. The text features a broad array of applications of computational methods to science and technology. This book would be an interesting supplement for the practicing engineers, scientists, and graduate students.

Published in London, UK

© 2018 IntechOpen
© Sean Hannon / iStock

IntechOpen

