

Proceedings of
the International
Conference

Numerical Methods for Scientific Computations and Advanced Applications

(NMSCAA'18)

In cooperation with
siam
and technically co-
sponsored by IEEE PS
Computer Society
Chapter (IEEE)

Editor
**Krassimir
Georgiev**

May 28 – May 31,
2018,
Hissar,
Bulgaria

**Sofia
2018**



Institute of Information and Communication Technologies
Bulgarian Academy of Sciences

Proceedings of the International Conference

**Numerical Methods for Scientific
Computations and Advanced Applications
(NMSCAA'18)**

May 28 – May 31, 2018, Hissarya, Bulgaria

Krassimir Georgiev (Editor)



Institute of Information and Communication Technologies
Bulgarian Academy of Sciences

In cooperation with



Technically co-sponsored by IEEE PS Computer Society Chapter

Sofia, 2018

Numerical Methods for Scientific Computations and Advanced Applications
(NMSCAA'18)

Proceedings of the International Conference

©2018 by Fastumprint

ISBN: 978-954-91700-7-8

Printed in Sofia, Bulgaria

PREFACE

This book contains papers presented during the International Conference on “Numerical Methods for Scientific Computations and Advanced Applications” (NMSCAA’18), May 28–May 31, 2018, Hissar, Bulgaria. The conference is organized by the Institute of Information and Communication Technologies, Bulgarian Academy of Sciences in cooperation with Society for Industrial and Applied Mathematics (SIAM) and technically co-sponsored by IEEE PS Computer Society Chapter (IEEE) .

The Conference Specific topics of interest are as follows: (a) Multiscale and multi-physics problems; (b) Robust preconditioning; (c) Monte Carlo methods; (d) Optimization and control systems; (e) Scalable parallel algorithms; (f) Advanced computing for innovations; (g) In silico investigations of biological molecules and complexes. The list of the plenary invited speakers includes:

- Owe Axelsson (Institute of Geonics, CAS, Ostrava, Czech Republic);
- Raytcho Lazarov (TA&MU, College Station, Texas, USA);
- Zahari Zlatev (Aarhus University, Roskilde, Denmark);
- Istvan Farago (Eotvos Lorand University, Budapest, Hungary);
- Radim Blaheta (Institute of Geonics, CAS, Ostrava, Czech Republic);
- Svetozar Margenov (IICT–BAS, Sofia, Bulgaria); and
- Ivan Dimov (IICT–BAS, Sofia, Bulgaria).

The Scientific Computing is one of the most prominent examples of a interdisciplinary area involving mathematics, computer science, engineering, physics, chemistry, medicine etc. The tools of Scientific Computing are usually based on mathematical models and corresponding computer codes that are used to perform virtual experiments to obtain new data or to better understand existing experimental results. Numerical Analysis is one of the crucial elements of Scientific Computing. It develops and analyzes numerical methods for discretization of continuous models and their subsequent solution, as well as for approximation of discrete data, such as: data interpolation and extrapolation, methods for solving linear and non-linear systems of algebraic equations (direct and iterative solution methods, preconditioning, multilevel and multigrid methods, etc.), methods for solving systems of ordinary and partial differential equations, methods for solving integral equations, and optimization problems.

Next to Numerical methods and the scientific computations are the Advanced Applications – the implementation of the developed numerical methods into computer codes and their customization for the numerous computing systems and for solving a number of real life problems.

Krassimir Georgiev

May 2018

Table of Contents

Part A: Short Communications/Extended abstracts	1
<i>G. Accaputo, P. Arbenz, P. Derlet</i> Solving Large-Scale Eigenvalue Problems in Amorphous Materials	3
<i>E. Atanassov, T. Gurov, M. Durchova, S. Ivanovska, A. Karavivanova</i> Study of Scalability and Energy Efficiency of QMC Algorithms on Hybrid HPC Systems	6
<i>O. Axelsson</i> Optimality Properties of a Square Block Matrix Preconditioner with Applications	10
<i>M. Beceanu, M. Lachaab</i> Fast Computation of Exact Solutions to the Heat and Stokes' Equations on the Half-Line Obtained by Fokas' Transform	14
<i>I. Blagoev, J. Sevova, K. Kolev</i> Artificial Neural Network Activation Function Optimization with Genetic Algorithms	16
<i>R. Blaheta, O. Axelsson, T. Luber, J. Kruzik, J. Stary</i> Preconditioners for Simulation of Flow in Rigid and Deformable Porous Media	20
<i>A. Cesmelioglu</i> A Monolithic Scheme for a Fluid-Poroelastic Structure Interaction Problem	24
<i>I. Dimov</i> Computational Nano-physics – Monte Carlo Approach	27
<i>B. Duan, R.D. Lazarov, J.H. Pasciak</i> Numerical Approximation of Fractional Spectral Elliptic Operators	28
<i>S.-E. Ekström, M. Neytcheva</i> Deflation Methods Made Possible	29
<i>I. Faragó, R. Horváth, M. Mincsovcics, R. Mosleh, F. Dorner</i> Reliable Numerical Models and Their Applications	32
<i>G. Gadzhev, K. Ganev</i> Vertical Structure of Atmospheric Composition Fields over Bulgaria	38
<i>I. Georgiev, I. Georgiev</i> Performance Analysis of Real-time Applications for Debugging Parametrization	42

<i>I. Georgieva, N. Miloshev</i> Computer Simulations of PM Concentrations Climate for Bulgaria	46
<i>S. Harizanov, N. Kosturski, R. Lazarov, S. Margenov, P. Marinov, Y. Vutov</i> Numerical Methods for Fractional-in-Space Diffusion Problems	50
<i>S. Harizanov, I. Lirkov, I. Georgiev, J. Stary, S. Zolotarev</i> Edge Detection of Radiographic Images through Phantom Blur Denoising	54
<i>Y. Hou, J. Dai, A. J. Niemi, X. Peng, J. He, N. Ilieva</i> Study of Non-Proline <i>cis</i> Peptide Planes in Different Protein Framings	56
<i>E. Lilkova, N. Ilieva, P. Petkov, L. Litov</i> Study of Human Interferon-Gamma Glycosilation by Molecular Dynamics Simulations	60
<i>K. Liolios, T. Makarios, A. Liolios, K. Georgiev, I. Georgiev</i> Monte Carlo Simulation for Seismic Analysis of Egnatia Highway Bridges in Northern Greece	64
<i>D. Slavchev, S. Margenov</i> Scalability Analysis of Solvers based on Hierarchical Compression of Dense Matrices and Gaussian Elimination	68
<i>B. Takács, R. Horváth, I. Faragó</i> A Non-Symmetric Model of Disease Propagation	72
<i>P. Tomov, I. Zankinski, M. Barova</i> Artificial Neural Networks Time Series Forecasting with Android Live Wallpaper Technology	76
<i>F. E. Uilhoorn</i> Pareto Optimal Solutions of Noise Statistics for Kalman Filtering Applied to State Estimation of Gas Dynamics	80
<i>Z. Zlatev, I. Dimov, I. Farago, K. Georgiev, A. Havasi</i> Implementation of the Three-times Repeated Richardson Extrapolation together with Explicit Runge-Kutta Methods	84
<i>O. Axelsson, S. Sysala</i> An adaptive Newton method for solving nonlinear partial differential equations	89
Part B: List of participants	93

Part A

Short Communications/Extended abstracts¹

¹Arranged alphabetically according to the family name of the first author.

Solving Large-Scale Eigenvalue Problems in Amorphous Materials

G. Accaputo, P. Arbenz, P. Derlet

In strongly amorphous materials, such as structural glasses, sound has anomalous dispersion properties which are characterized by the Ioffe–Regel limit, and the universal phenomenon of the Boson peak. Molecular dynamics simulations are able to produce structural glasses in which the stable position of each atom is precisely known. Given such a computer generated atomic configuration, the corresponding vibrational properties can be investigated by solving an eigenvalue problem involving the Hessian of the potential energy landscape associated with material cohesion.

The Hessian obtained from a large scale molecular dynamics simulation of a model metallic structural glass involving millions to billions of atoms is to be constructed and *partially* diagonalized to obtain the significant portion of the long-wavelength vibrational eigenmodes [3]. These eigenmodes will be analyzed to investigate the relationship between atomic-scale structure and the onset of the Boson peak regime. In order to calculate the vibrational frequencies, we have to solve the real symmetric eigenvalue problem

$$\mathbf{H}\mathbf{q} = \lambda\mathbf{q}, \tag{1}$$

where \mathbf{H} is the given *Cartesian Hessian* that contains the second derivatives of the total electronic energy with respect to nuclear Cartesian coordinates. The presented simulation involves approximately $1.5 \cdot 10^6$ atoms leading to a matrix \mathbf{H} of size about $4.5 \cdot 10^6$. We are interested in the 100 to 1000 eigenvalues λ_k and associated eigenvectors \mathbf{q}_k of (1) that are closest to the Boson peak.

The shift-and-invert Lanczos (SI-Lanczos) algorithm is the method of choice for computing interior eigenvalues and corresponding eigenvectors of a symmetric or Hermitian matrix \mathbf{H} close to some target τ . However, the SI-Lanczos algorithm needs the factorization of $\mathbf{H} - \tau\mathbf{I}$ which is not feasible here for its excessive memory requirements. For such cases, the Jacobi–Davidson methods have been developed [11]. To be efficient, they however need an effective preconditioner to solve the so-called correction equation, which usually entails its factorization. In an earlier study [8], we were not able to identify such preconditioners for (1).

In this work we investigate a technique, known as *spectral filtering*, for solving eigenvalue problems that obviates factorizations altogether [10, 7] Spectral filtering is combined in practice with Krylov space methods [4, 9] and subspace iteration [13, 5]. In order for the technique to be applicable the extremal eigenvalues $\lambda_{\min}(\mathbf{H})$ and $\lambda_{\max}(\mathbf{H})$ of \mathbf{H} , or, at least, some decent bounds must be known. To compute the eigenvalues in the interval $[\xi, \eta] \subset [\lambda_{\min}, \lambda_{\max}]$, a polynomial $\rho \in \mathbb{P}_d$ is constructed such that The desired polynomial ρ could be an approximation of the characteristic function $\chi_{[\xi, \eta]}$ associated with the interval $[\xi, \eta]$. If $\rho(\mathbf{H})$ multiplies a vector, (most) of the unwanted eigenvector components are suppressed. Therefore, ρ is called a *polynomial filter*. The degree of ρ depends on the width of the interval $[\xi, \eta]$, on the width

ε of the margins, and the strength of the filter. The degree increases if $\eta - \xi$ and/or ε shrink. In our experiments we use polynomial degrees d as high as $\mathcal{O}(1000)$. A consequence is that increasing parallelism by slicing the interval $[\xi, \eta]$ is not scalable. Interval slicing may however be necessary for memory reasons.

After a brief review of the technique of polynomial filtering we suggest a filter that should be useful for filtering in connection with the subspace iteration method. We investigate the filters' properties by means of some synthetic eigenvalue problems. Based on these findings we investigate a model of amorphous solid consisting of 1'372'000 atoms corresponding to a \mathbf{H} of size 4'116'000.

Our eigensolvers is implemented with the Trilinos software framework [1]. Trilinos [6, 12] is a collection of open-source software libraries, called *packages*, for the development of scientific applications. *Anasazi* [2] is a package that offers a collection of algorithms for solving large-scale eigenvalue problems. We employ Anasazi's block Krylov–Schur eigensolver with thick restarts. The subspace iteration we implemented it ourselves, based on Trilinos *Epetra* data structures. The large scale computations have been carried out on the Euler cluster of ETH Zürich².

References

- [1] G. Accaputo. Solving large scale eigenvalue problems in amorphous materials. Master's thesis, ETH Zurich, Computer Science Department, September 2017. doi:10.3929/ethz-b-000221499.
- [2] C. G. Baker, U. L. Hetmaniuk, R. B. Lehoucq, and H. K. Thornquist. Anasazi software for the numerical solution of large-scale eigenvalue problems. *ACM Trans. Math. Softw.*, 36(3):1–23, 2009.
- [3] P. M. Derlet, R. Maaß, and J. F. Löffler. The Boson peak of model glass systems and its relation to atomic structure. *Eur. Phys. J. B*, 85(5):1–20, 2012.
- [4] H.-R. Fang and Y. Saad. A filtered Lanczos procedure for extreme and interior eigenvalue problems. *SIAM J. Sci. Comput.*, 34(4):A2220–A2246, 2012.
- [5] M. Galgon, L. Krämer, B. Lang, A. Alvermann, H. Fehske, A. Pieper, G. Hager, M. Kreutzer, F. Shahzad, G. Wellein, A. Basermann, M. Röhrig-Zöllner, and J. Thies. Improved coefficients for polynomial filtering in ESSEX. In T. Sakurai, S.-L. Zhang, T. Imamura, Y. Yamamoto, Y. Kuramashi, and T. Hoshi, editors, *Eigenvalue Problems: Algorithms, Software and Applications in Petascale Computing*, pages 63–79. Springer, 2017.
- [6] M. A. Heroux et al. An overview of the Trilinos project. *ACM Trans. Math. Softw.*, 31(3):397–423, 2005.

²<https://scicomp.ethz.ch/wiki/Euler>

- [7] L. O. Jay, H. Kim, Y. Saad, and J. R. Chelikowsky. Electronic structure calculations for plane-wave codes without diagonalization. *Comput. Phys. Comm.*, 118(1):21 – 30, 1999.
- [8] S. Schaffner. Using Trilinos to solve large scale eigenvalue problems in amorphous materials. Master’s thesis, ETH Zurich, Computer Science Department, April 2015.
- [9] G. Schofield, J. R. Chelikowsky, and Y. Saad. A spectrum slicing method for the Kohn–Sham problem. *Comput. Phys. Comm.*, 183(3):497–505, 2012.
- [10] R. N. Silver, H. Röder, A. F. Voter, and J. D. Kress. Kernel polynomial approximations for densities of states and spectral functions. *J. Comput. Phys.*, 124(1):115–130, 1996.
- [11] G. L. G. Sleijpen and H. A. van der Vorst. A Jacobi–Davidson iteration method for linear eigenvalue problems. *SIAM J. Matrix Anal. Appl.*, 17(2):401–425, 1996.
- [12] The Trilinos Project Home Page. <http://trilinos.org/>.
- [13] Y. Zhou, Y. Saad, M. L. Tiago, and J. R. Chelikowsky. Self-consistent field calculations using Chebyshev-filtered subspace iteration. *J. Comput. Phys.*, 219(1):172–184, 2006.

Study of Scalability and Energy Efficiency of QMC Algorithms on Hybrid HPC Systems

E. Atanassov, T. Gurov, M. Durchova, S. Ivanovska,
A. Karaivanova

1 Introduction

In this work we study the scalability issues and the energy efficiency of some of the typical QMC algorithms for hybrid HPC systems. Depending on the features of the algorithms and the selected sequences, we investigate the optimal setup of MPI processes and OpenMP threads from point of view of speed and energy efficiency. The positive impact from using the vector instructions of the Intel Xeon Phi accelerators is demonstrated compared with CPUs and with automatic vectorization by the compiler. This research is motivated by our experience with the supercomputer Avitohol at IICT-BAS [6] which total performance is 412 TFlops, but 90% comes from the accelerators.

The Xeon Phi coprocessors combine efficient vector floating point computations with familiar operational and development environment. It is already known that to obtain good performance on MIC architecture, algorithms have to be vectorized to exploit the vector engine of the processor. On such specialised equipment like the Xeon Phi, the importance of memory accesses for the overall performance increases due to the presence of a large number of computational cores. The different parallelisation models entail different trade-off considerations between use of more memory or making more computations. The hybrid OpenMP+MPI programming is more complicated, but potentially is the most advantageous, especially when high numbers of cores and servers are utilized [8]. The increased weight of energy cost justifies our focus on developing energy efficient algorithms. In this work, we present our parallelisation strategies and representative case studies (low discrepancy sequence generation and solving multidimensional integrals). Numerical and timing results are shown and discussed.

2 Scalability and Timing Results

During our investigation, we compared generation codes for the Halton [1] and Sobol sequences [5], obtained via using the auto-vectorization features of the Intel compiler and hand-tuned vectorized codes, as it is shown on Fig. 1. For the Sobol sequence the auto-vectorization resulted in approximately two times better performance, while for the Halton sequence there was a negligible difference. This can be explained by the higher complexity of the Halton sequences definition. Nevertheless, in both cases the vectorization by hand lead to several times better performance. In the case of the Halton sequences, this was achieved via substantial reorganization of the

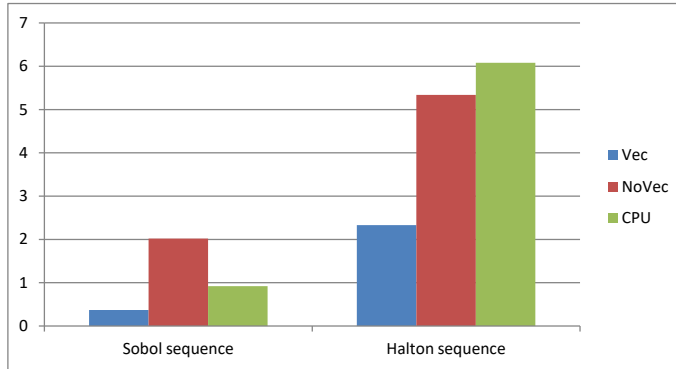


Figure 1: Timing results for Sobol and Halton sequences using three porting methods

operations. The hand-tuned Sobol generation code was even faster than the standard Intel implementation of the Mersenne twister pseudorandom number generator in the MKL [7].

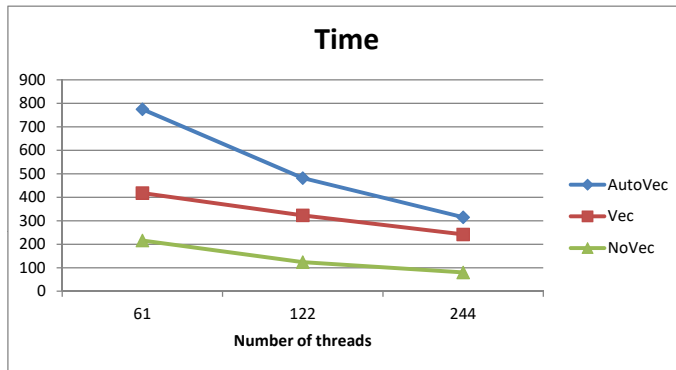


Figure 2: Scalability results (Sobol sequence) for three variations of the algorithm (auto-vectorized, vectorized by hand, without vectorization)

The Intel Xeon Phi coprocessor has 244 hardware threads, which means that the available 61 physical cores can be used with up to 4x hyper-threading. Our results show that the use of 244 threads can give substantial performance benefit versus the use of only 61 threads. That is why we recommend users to utilize as many threads as possible when using our generation routines. On Fig. 2 one can see the timing results of the various ways of generating the Sobol sequences (the results for Halton are similar).

3 Study of Energy Efficiency

The energy use of the generators was investigated, while also varying the number of hardware threads of the Intel Xeon Phi co-processor. One can see that the use of more hardware threads leads to higher energy usage (see [2], [3], [4]). That is why there is a thread-off between the time-to-completion and price of energy that will determine which number of threads is optimal. It appears that in many cases 122 threads will be the optimal number to use (see Fig. 3 and Fig. 4).

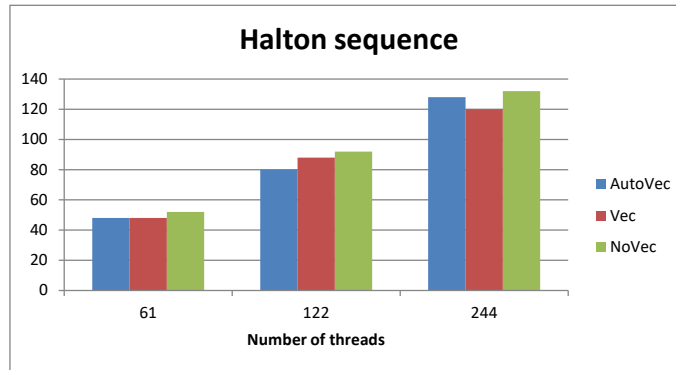


Figure 3: Halton sequence energy usage of the algorithm with three optimizations (auto-vectorized, vectorized by hand, without vectorization)

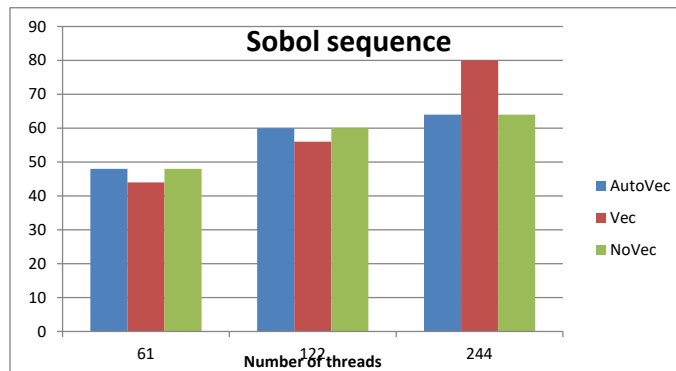


Figure 4: Sobol sequence energy usage of the algorithm with three optimizations (auto-vectorized, vectorized by hand, without vectorization)

4 Conclusion

We developed and compared vectorized variants of generation routines for the Sobol and Halton low-discrepancy sequences for the Intel Xeon Phi accelerators implementation. With necessary expertise (and time) the researchers might develop highly optimized algorithms which maximize the benefits of using the accelerators. In the case of use of hyperthreading we observe sizeable benefit from using the larger number of threads, which should be weighted against increase energy consumption.

Acknowledgments

This work was supported by the National Science Fund of Bulgaria under Grant DFNI-I02/8 and by the European Commission under H2020 project VI-SEEM (contract number 675121).

References

- [1] E. Atanassov, M. Durchova, Generation of the Scrambled Halton Sequence Using Accelerators, MIPRO 2013, Proceedings of the 36th International Convention, IEEE, pp. 197–201, ISSN: 1847–3946.
- [2] E. Atanassov, T. Gurov, A. Karaivanova, Energy Aware Performance Study for a Class of Computationally Intensive Monte Carlo Algorithms, J. Computers & Mathematics with Applications, Volume 70, Issue 11, 2015, pp. 2719–2725
- [3] C. Bekas, A. Curioni, A new energy aware performance metric, Springer, Comput. Sci. Res Dev (2010) 25: 187–195, DOI 10.1007/s00450-010- 0119-z.
- [4] J. Demmel, A. Gearhart, B. Lipshitz, O. Schwartz, Perfect Strong Scaling Using No Additional Energy. Proc. of IEEE 27th IPDPS13. IEEE Computer Society, 2013
- [5] I. Sobol, D. Asotsky, A. Kreinin, S. Kucherenko (2011). Construction and Comparison of High-Dimensional Sobol Generators. Wilmott Journal Nov: 64–79.
- [6] High-performance computing system — Avitohol, <http://www.hpc.acad.bg/system-1/>
- [7] Intel Math Kernel Library (MKL), <http://software.intel.com/en-us/articles/intel-math-kernel-library-documentation>
- [8] MPI Standard: <http://www.mcs.anl.gov/research/projects/mpi/>

Optimality Properties of a Square Block Matrix Preconditioner with Applications

O. Axelsson

Dedicated to Krassimir Georgiev, Deputy Director, Institute of Information and Communication Technologies, BAS, Sofia, Bulgaria, as a reminder of a very long lasting friendship.

Two-by-two block matrices with block rows A , $-B^\top$, and B , A , arise in a number of important applications, such as when solving complex valued matrix systems in real valued form, in optimal control problems for PDEs with various kind of state stationary or time dependent equations, such as Poisson, convection diffusion, Stokes [2], and Maxwell [4] and in wave propagation and structural dynamics and many more. Since the discretized problems have a large scale, iterative solution methods must be used and then preferably with a preconditioner with optimal properties. The first part of this paper presents such a preconditioner which depends on a parameter. Assuming that A and the sum of the B -matrices are symmetric and positive semidefinite and A and B have disjoint nullspaces, the optimal value of the parameter is derived for two versions of an iterative refinement method and for the application of the preconditioner with Chebychev or conjugate gradient methods as acceleration methods. The resulting eigenvalue bounds are very tight and the rate of convergence factor for each iteration step (after rounding), takes a value between 0.333 and 0.172 for the different methods. This holds uniformly with respect to all classes of the above problems and for various problem and method parameters. The iterative refinement methods including the Chebyshev method do not require computations of inner products as the CG type methods do, which can save much communication overhead and elapsed computer times on massively parallel computer platforms. The preconditioner requires two solutions of a matrix system with a fixed linear combination of the block row matrices for which efficient iterative solution methods exist. The methods outperform other published methods for the above classes of problems. In the second part of the paper, application of the methods for an optimal control problem with a PDE of parabolic type and an eddy current electromagnetic problem are presented. In the electromagnetic problem the electrical field is used as control of the magnetic solution to take a desired shape.

Let the preconditioner be defined by the following parameter version of the PRESB method method, see [2], with matrices of order $n \times n$, where previously we have chosen $\alpha = 1$,

$$C = \begin{bmatrix} A & -B^\top \\ B & \alpha^2 A + \alpha(B + B^\top) \end{bmatrix}. \quad (1)$$

A solution method based on Schur complements involve actions of A^{-1} , which must be avoided as A can be singular. We show now how this can be done. Consider then

$$\mathcal{C} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} f \\ g \end{bmatrix} \quad (2)$$

where we multiply the first equation by α and add the second equation, to form

$$\begin{cases} (\alpha A + B)(x + \alpha y) &= \alpha f \\ B(x + \alpha y) + \alpha(\alpha A + B^T)y &= g \end{cases}.$$

This shows that, besides some vector operators and a matrix vector multiplication with B , its solution requires just one solution with $H = \alpha A + B$ and one with H^T . To find eigenvalues of $\mathcal{C}^{-1}\mathcal{A}$, consider now the generalized eigenvalue problem,

$$\lambda \mathcal{C} \begin{bmatrix} x \\ y \end{bmatrix} = \mathcal{A} \begin{bmatrix} x \\ y \end{bmatrix}, \quad \|x\| + \|y\| \neq 0. \quad (3)$$

This can be rewritten as

$$(1 - \lambda) \mathcal{C} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 0 \\ ((\alpha^2 - 1)A + \alpha(B + B^T))y \end{bmatrix}.$$

Since

$$(1 - \lambda)(Ax - B^T y) = 0$$

it follows that $\lambda = 1$ if $Ax \neq B^T y$. Hence the multiplicity of the unit eigenvalue is at least n . If $Ax = B^T y$, it follows that

$$\frac{1}{1 + \alpha^2} \leq \lambda \leq \frac{1}{\alpha^2} \quad \text{for } \alpha \leq 1, \quad \frac{1}{1 + \alpha^2} \leq \lambda \leq 1 \quad \text{for } \alpha \geq 1.$$

Hence the spectral condition number

$$\kappa(\mathcal{C}^{-1}\mathcal{A}) \leq 1 + \frac{1}{\alpha^2} \quad \text{for } \alpha \leq 1 \quad \text{and} \quad \kappa(\mathcal{C}^{-1}\mathcal{A}) = 1 + \alpha^2 \quad \text{for } \alpha \geq 1,$$

which are both minimized and takes value 2 for $\alpha = 1$. It follows that the spectral condition number bound for the preconditioned PRESB method where $\alpha = 1$, see [2, 3], can not be improved by use of such perturbations. It is further seen that the preconditioned matrix has the block form

$$\mathcal{C}^{-1}\mathcal{A} = \begin{bmatrix} I & F \\ 0 & I - E \end{bmatrix},$$

where $E = H^{-T}AH^{-1}G$, $G = (\alpha^2 - 1)A + \alpha(B + B^T)$. Since $H^{-T}AH^{-1}G$ has a complete eigenvector space it follows that E is a normal matrix and hence, that the preconditioned conjugate gradient method in exact arithmetic gives the minimal polynomial, i.e. best polynomial approximation of the solution. Hence no breakdown of the iterations can occur.

Consider now a defect-correction method with an iterative refinement parameter $\tau > 0$, to solve a linear system $\mathcal{A}x = b$, with by use of the preconditioner,

$$\mathcal{C}(x^{k+1} - x^k) = \tau(b - \mathcal{A}x^k), \quad k = 0, 1, \dots, x^0 \text{ given.}$$

Hence

$$e^{k+1} = (1 - \tau)e^k + \tau\mathcal{C}^{-1}Re^k, \quad k = 0, 1, \dots.$$

where $e^k = \hat{x} - x^k$ and $R = \mathcal{C} - \mathcal{A}$.

The eigenvalues $\mu(\mathcal{C}^{-1}R)$ are contained in the intervals,

$$\frac{\alpha^2}{\alpha^2 + 1} \geq \mu(\mathcal{C}^{-1}R) \geq \frac{\alpha^2 - 1}{\alpha^2}, \quad \text{for } \alpha \leq 1,$$

so the convergence factor is minimized for $\tau = \frac{2\alpha^2(\alpha^2+1)}{2\alpha^2+1}$ and the rate of convergence factor becomes

$$1 - \frac{\tau}{\alpha^2 + 1} = 1 - \frac{2\alpha^2}{2\alpha^2 + 1} = \frac{1}{2\alpha^2 + 1},$$

which, since $\alpha \leq 1$, is minimized for $\alpha = 1$. The optimal value of τ is then $\tau = \frac{4}{3}$ and $\|e^{k+1}\| \leq \frac{1}{3}\|e^k\|$, $k = 0, 1, \dots$.

The iterative refinement method can be seen as an Euler forward time-stepping method with time-step τ to solve

$$\mathcal{C} \frac{dx(t)}{dt} = b - \mathcal{A}x(t), \quad t > 0, \quad x(0) = x_0.$$

Clearly we can repeat this method using two or more different time-steps. For two consecutive time-steps with $\tau_1 > 0$, $\tau_2 > 0$, we get

$$\begin{cases} \mathcal{C}x^{k+1/2} &= (1 - \tau_1)\mathcal{C}x^k + \tau_1 Rx^k \\ \mathcal{C}x^{k+1} &= (1 - \tau_2)\mathcal{C}x^{k+1/2} + \tau_2 Rx^{k+1/2}. \end{cases}$$

Hence

$$x^{k+1} = (1 - \tau_2)(1 - \tau_1)x^k + ((1 - \tau_2)\tau_1 + \tau_2(1 - \tau_1))\mathcal{C}^{-1}Rx^k + \tau_1\tau_2(\mathcal{C}^{-1}R)^2x^k$$

or

$$x^{k+1} = Q_2x^k, \quad k = 0, 1, \dots$$

where

$$\begin{aligned} Q_2 &= (1 - \tau_2)(1 - \tau_1)I + (\tau_1 + \tau_2 - 2\tau_1\tau_2)\mathcal{C}^{-1}R + \tau_1\tau_2(\mathcal{C}^{-1}R)^2 = \\ &= I - (\tau_1 + \tau_2)(I - \mathcal{C}^{-1}R) + \tau_1\tau_2(I - \mathcal{C}^{-1}R)^2. \end{aligned}$$

As is wellknown, this is minimized if Q_2 equals the normalized second degree Chebyshev polynomial

$$Q_k = T_k \left(\frac{a+b-2\xi}{b-a} \right) / T_k \left(\frac{b+a}{b-a} \right), \quad (k = 2)$$

where a, b are the eigenvalue bounds and

$$T_{k+1}(\xi) = 2\xi T_k(\xi) - T_{k-1}(\xi), \quad k = 1, 2, \dots, \quad T_0(\xi) = 1, \quad T_1(\xi) = \xi.$$

Hence $T_2(\xi) = 2\xi^2 - 1$. In our problem, the eigenvalue bounds are $a = \frac{1}{\alpha^2+1}$, $b = \frac{1}{\alpha^2}$, i.e., $a = \frac{1}{2}$, $b = 1$ for the optimal value $\alpha = 1$, so

$$Q_2 = \frac{2 \left(\frac{a+b-2\xi}{b-a} \right)^2 - 1}{2 \left(\frac{a+b}{b-a} \right)^2 - 1} = \frac{17 - 48\xi + 32\xi^2}{17},$$

that is, $\tau_1 + \tau_2 = \frac{48}{17}$, $\tau_1\tau_2 = \frac{32}{17}$, and $\tau_i = \frac{4}{17}(6 \pm \sqrt{2})$, $i = 1, 2$. Further

$$\|Q_2\| = \frac{1}{T_2 \left(\frac{b+a}{b-a} \right)} = \frac{1}{2 \cdot 9 - 1} = \frac{1}{17}.$$

Since this method involves two iteration steps it corresponds to an average convergence factor $\frac{1}{\sqrt{17}} \approx 0,241$ per iteration step. Extending the method with more iterative refinement steps with k different time steps corresponds to the recursion

$$x^{k+1} = Q_k x^k, \quad k = 0, 1, \dots$$

where the optimal matrix polynomial Q_k is the Chebyshev polynomial.

In this way we approach the rate of convergence factor $\frac{\sqrt{2}-1}{\sqrt{2}+1} = \frac{1}{(\sqrt{2}+1)^2} \simeq 0,173$ which holds for the preconditioned Chebyshev iteration method.

References

- [1] O. Axelsson, M. Neytcheva, B. Ahmad. A comparison of iterative methods to solve complex valued linear algebraic systems. *Numerical Algorithms*, 66(2014) 811–841.
- [2] O. Axelsson, S. Farouq, M. Neytcheva. Comparison of preconditioned Krylov subspace iteration methods for PDE-constrained optimization problems. Stokes control. *Numerical Algorithms*, 74(2017) 19–37.
- [3] Z.-Z. Liang, O. Axelsson, M. Neytcheva. A robust structured preconditioner for the time-harmonic parabolic optimal control problem. Accepted by *Numerical Algorithms*, 2018.
- [4] O. Axelsson, D. Lukáš. Preconditioning methods for eddy current optimally controlled time-harmonic electromagnetic problems. *J. Numerical Mathematics*, (2018), to appear.

Fast Computation of Exact Solutions to the Heat and Stokes' Equations on the Half-Line Obtained by Fokas' Transform

M. Beceanu, M. Lachaab

A new method has been recently introduced by Fokas (Fokas, 2002 [3]; Fokas et al., 2009 [4]) for solving a large class of PDEs. This method, however, raises a problem: the solutions of the PDEs are given as contour integrals on an unbounded contour in the complex plane that need to be evaluated numerically. To evaluate these integrals, Flyer–Fokas (2008, [2]), and Papatheodorou–Kandili (2009, [5]) deformed and parametrized the contour of integration and used the simple trapezoid rule without analyzing the error.

In this paper, we obtain exact expressions for the solutions of the heat equation and Stokes' equation of the first kind in terms of elementary functions, the imaginary error function, and the incomplete Airy function. For the heat equation with zero initial condition and sine boundary condition, the solution is given by:

$$q(x, t) = \frac{e^{-\frac{x^2}{4t}}}{2} \operatorname{Re} \left[e^{-z_1^2} (i + \operatorname{erfi}(-z_1)) + e^{-z_2^2} (i + \operatorname{erfi}(-z_2)) \right],$$

where

$$z_1 = -\sqrt{\frac{\lambda t}{2}} + i \left(\frac{x}{2\sqrt{t}} - \sqrt{\frac{\lambda t}{2}} \right), \quad z_2 = \sqrt{\frac{\lambda t}{2}} + i \left(\frac{x}{2\sqrt{t}} + \sqrt{\frac{\lambda t}{2}} \right).$$

And for general boundary data $g_0(t)$, the solution is given by:

$$q(x, t) = \frac{e^{-\frac{x^2}{4t}}}{\sqrt{2\pi}} \int_0^\infty \operatorname{Re} \left[e^{-z_1^2} (i + \operatorname{erfi}(-z_1)) + e^{-z_2^2} (i + \operatorname{erfi}(-z_2)) \right] \widehat{g}_0^s(\lambda) d\lambda,$$

where \widehat{g}_0^s is the sine transform of g_0 :

$$\widehat{g}_0^s(\lambda) = \frac{2}{\pi} \int_0^\infty g_0(t) \sin(\lambda t) dt.$$

For the Stokes' equation, the solution is given by:

$$\begin{aligned} q(x, t) &= \frac{1}{2} \left(\sin(\lambda t + \lambda^{\frac{1}{3}} x) + e^{-(\frac{\sqrt{3}}{2} + \frac{i}{2})\lambda^{\frac{1}{3}} x} \sin(\lambda t) \right) \\ &+ \frac{1}{4} \left(\sum_{j=1}^2 e^{i(\frac{\tilde{\alpha}_j^3}{3} + \tilde{x}\tilde{\alpha}_j)} \left(i - \frac{f_{\tilde{x}}(\tilde{\alpha}_j)}{C_{\tilde{x}}} \right) \right. \\ &\left. - \sum_{j=3}^4 e^{i(\frac{\tilde{\alpha}_j^3}{3} + \tilde{x}\tilde{\alpha}_j)} \frac{f_{\tilde{x}}(\tilde{\alpha}_j)}{C_{\tilde{x}}} + \sum_{j=5}^6 e^{i(\frac{\tilde{\alpha}_j^3}{3} + \tilde{x}\tilde{\alpha}_j)} \left(2i - \frac{f_{\tilde{x}}(\tilde{\alpha}_j)}{C_{\tilde{x}}} \right) \right) \end{aligned}$$

where: $f_x(k) = \int_k^{+i\infty} e^{-i(\frac{t^3}{3}+tx)} dt$, $C_x = \int_0^{+\infty} e^{-\frac{t^3}{3}+tx} dt$, $\tilde{x} = (\frac{1}{3t})^{\frac{1}{3}}x$, and $\tilde{\alpha}_j = (3t)^{\frac{1}{3}}\alpha_j$.

The above solutions lend themselves well to numerical computations, since there exist fast and highly accurate methods for computing the imaginary error function and the incomplete Airy function. For example, the imaginary error function is a standard function in MATLAB, where it is estimated with 10^{-20} accuracy by means of Padé approximants (Cody, 1969 [1]). Also, the above solutions extend to the lateral boundaries without convergence issues, allow for an easy analysis of the estimation error, and are much faster than those obtained by other methods.

In addition, we derive an asymptotic expansion for the solution of the heat equation with precise bounds for the error term, which allows one to compute the solution with arbitrarily high precision. The solution to the heat equation admits the following asymptotic expansion:

$$q(x, t) = \begin{cases} e^{-x\sqrt{\frac{\lambda}{2}}} \sin\left(\lambda t - x\sqrt{\frac{\lambda}{2}}\right) + u(x, t), & x \ll t\sqrt{2\lambda} \\ u(x, t), & x \gg t\sqrt{2\lambda}, \end{cases}$$

where

$$u(x, t) = -\frac{e^{-\frac{x^2}{4t}}}{2\sqrt{\pi}} \sum_{n=0}^{N-1} \operatorname{Re}\left(\frac{(2n-1)!!}{2^{n+1}z_1^{2n+1}} + \frac{(2n-1)!!}{2^{n+1}z_2^{2n+1}}\right) + Q'_{2N}$$

and the error Q'_{2N} is bounded by

$$|Q'_{2N}| < \frac{2e^{-\frac{x^2}{4t}}\sqrt{t}(2N-1)!!}{\max\left(\lambda t, \frac{x^2}{4t}\right)^N |x - t\sqrt{2\lambda}|\sqrt{\pi}}.$$

References

- [1] Cody, W.J., (1969), “Rational Chebyshev approximations for the error function”, *Mathematics of Computation*, Vol 23, 107, 631-637.
- [2] Flyer, N., Fokas, A.S., (2008), “A hybrid analytical numerical method for solving evolution partial differential equations in the half-line”, *Proceedings of the Royal Society A*, 464, pp. 1823–1849.
- [3] Fokas, A.S., (2002), “A new transform method for evolution PDEs”, *IMA J. Appl. Math*, 67, 559 – 590.
- [4] Fokas, A.S., et al., (2009), “A semi analytical numerical method for solving evolution and elliptic partial differential equation”, *J. of Comp. and Applied Math*, 227, 59 – 74.
- [5] Papatheodorou, T.S., Kandili A.N., (2009), “Novel numerical techniques based on Fokas transforms for the solution of initial boundary value problems”, *J. of Comp. and Applied Math*, 227, 75 – 82.

Artificial Neural Network Activation Function Optimization with Genetic Algorithms

I. Blagoev, J. Sevova, K. Kolev

Abstract

Artificial Neural Networks (ANNs) are widely used in the last few decades. They have application in many different areas as financial forecasting [1, 7, 11], board and puzzle games [4, 6, 8], image processing, object classification [2], computer networks [9, 10] and many others. The most popular ANNs are represented as directed weighted graph in which signals are traveling from the input to the output. Each node (neuron) in these models have an activation function according which the node signal is emitted. The most popular activation functions are the sigmoid function and the hyperbolic tangent function [12, 5]. In this study evolutionary algorithms are employed in order to search for activation function alternatives. The function itself is represented as mathematical expression as the used in the genetic programming (GP).

Keywords: artificial neural networks, activation function, genetic algorithms, genetic programming

1 Introduction

The main purpose of neuron's the activation function is to limit the strength of the emitted signal. The most common way of input signal collection is by sum of the signals multiplied by weight of the connections (Eq. 1).

$$y_j = \sum_{i=1} x_i * w_{ij} \quad (1)$$

Collecting the signals in such way is very dependant of the size in the previous layer of neurons. Also the values of the weights varies a lot and the result of the Eq. 1 can reach high negative or positive values. Both problems are solved with normalization by usage of a neuron activation function.

$$z_j = \frac{1}{1 + e^{-y_j}} \quad (2)$$

$$z_j = \frac{e^{2y_j} - 1}{e^{2y_j} + 1} \quad (3)$$

The most used functions are the sigmoid function (Eq. 2) and the hyperbolic tangent function (Eq. 3). List of less common used activation functions can be found at [14].

2 Activation Function Optimization

The idea used in Radial Basis Function (RBF) ANNs for parameters optimization with Genetic Algorithms (GAs) [3] can be extended to optimization of the activation function itself. The goal in such optimization is to find activation function expression which will reduce ANN training time without lost of ANN operation accuracy. Implementing ANN solutions with the Encog Machine Learning Framework allows usage of alternative activation functions [12]. The expression of the alternative activation function can be easily represented with mXparser math expression library. The expression is stored and evaluated as text string. Activation function expressions as strings perfectly fit to the concept of the GP. Each string is represented as string chromosome and Apache Genetic Algorithms Framework is used for the chromosome evaluations. Initial GA population consists of valid randomly generated mathematical expressions. Single cut crossover is used and the random cut point is always on mathematical operator. As mutation random replacement with randomly selected mathematical operator is used. For fitness evaluate each expression is loaded in ANN neurons and total ANN efficiency is calculated. The total error for the neural network is used as fitness value.

3 Conclusions

Investigations in the direction of searching for alternative activation functions can achieve interesting results, which may change the way in which ANNs are used. As further research it can be interesting the experiments with activation function to proceed in combination with permutational algorithms as it was described in [13].

Acknowledgements

This work was supported by private funding of Velbazhd Software LLC.

References

- [1] Atanasova, T., Barova, M., Balabanov, T., *Use of Neural Models for Analysis of Time Series in Big Data*, Publishing complex of "Vasil Levski" National Military University, ISSN 1314–1937, 193–198, 2016.
- [2] Atanasova, T., Atanasov, J., *Business Processes Traceability in SME by Barcode System*, Proceedings of the International Scientific Conference, UNITECH'16, Gabrovo, Bulgaria, ISSN 1313–230X, 207–212, 2016.
- [3] Awad, M., *Optimization RBFNNs Parameters using Genetic Algorithms: Applied on Function Approximation*, International Journal of Computer Science and Security, vol. 4, issue 3, ISSN 1985–1553, 295–307, 2010.

- [4] Balabanov, T., Genova, K., *Distributed System for Artificial Neural Networks Training Based on Mobile Devices*, Proceedings of the International Conference Automatics and Informatics, Sofia, Bulgaria, Federation of the Scientific Engineering Unions John Atanasoff Society of Automatics and Informatics, ISSN 1313–1850, 49–52, 2016.
- [5] Balabanov, T., Keremedchiev, D., Goranov, I., *Web Distributed Computing For Evolutionary Training Of Artificial Neural Networks*, International Conference InfoTech, Varna — St. St. Constantine and Elena resort, Bulgaria, Publishing House of Technical University — Sofia, ISSN 1314–1023, 210–216, 2016.
- [6] Balabanov, T., Zankinski, I., Barova, M., *Strategy for Individuals Distribution by Incident Nodes Participation in Star Topology of Distributed Evolutionary Algorithms*, Cybernetics and Information Technologies, Institute of Information and Communication Technologies — BAS, vol. 16, no. 1, ISSN 1311–9702, 80–88, 2016.
- [7] Balabanov, T., Zankinski, I., Dobrinkova, N., *Time Series Prediction by Artificial Neural Networks and Differential Evolution in Distributed Environment*. Proceedings of the International Conference on Large-Scale Scientific Computing, Sozopol, Bulgaria, Lecture Notes in Computer Science, Springer, vol. 7116, no. 1, ISBN 978-3-642-29842-4, 198–205, 2011.
- [8] Keremedchiev, D., Barova, M., Tomov, P., *Mobile Application as Distributed Computing System for Artificial Neural Networks Training Used in Perfect Information Games*, Proceedings of the International Scientific Conference, UNITECH'16, Gabrovo, Bulgaria, ISSN 1313–230X, 389–393, 2016.
- [9] Tashev, T., Marinov, M., Monov, V., Tasheva, R., *Modeling of the MiMa-algorithm for crossbar switch by means of Generalized Nets*, Proceedings of the 2016 IEEE 8th International Conference on Intelligent Systems (IS), Sofia, Bulgaria, ISBN 978-1-5090-1354-8, 593–598, 2016.
- [10] Tashev, T., Monov, V., *Modeling of the hotspot load traffic for crossbar switch node by means of generalized nets*, Proceedings of the 6-th International IEEE Conference Intelligent Systems IS'12, Sofia, Bulgaria, vol. 2, 187–191, 2012.
- [11] Tomov, P., Monov, V., *Artificial Neural Networks and Differential Evolution Used for Time Series Forecasting in Distributed Environment*, Proceedings of the International Conference Automatics and Informatics, Sofia, Bulgaria, ISSN 1313–1850, 129–132, 2016.
- [12] Zankinski, I., Tomov, P., Balabanov, T., *Alternative Activation Function Derivative in Artificial Neural Networks*, 25th Symposium with International Participation — Control of Energy, Industrial and Ecological Systems, Bankia, Bulgaria, John Atanasoff Union of Automation and Informatics, ISSN 1313–2237, 79–81, 2017.

- [13] Zankinski, I., Stoilov, T., *Effect of the Neuron Permutation Problem on Training Artificial Neural Networks with Genetic Algorithms in Distributed Computing*, Proceedings of the 24th International Symposium Management of Energy, Industrial and Environmental Systems, ISSN 1313-2237, Bankya, Bulgaria, 53-55, 2016.
- [14] Wikipedia, *Activation Function*,
https://en.wikipedia.org/wiki/Activation_function

Preconditioners for Simulation of Flow in Rigid and Deformable Porous Media

R. Blaheta, O. Axelsson, T. Luber, J. Kruzik, J. Stary

1 Introduction and Considered Problems

The fluid flow in porous media appears in many applications dealing with geomaterials, biomaterials etc. In the case of fully saturated rigid porous media, the flow problem can be formulated by using the pressure in fluid $p = p(x, t)$, the fluid velocity $v = v(x, t)$ and the following two equations

$$\begin{aligned} k^{-1}v + \nabla p &= 0, \\ \operatorname{div}(v) + c_{pp} \frac{\partial}{\partial t} p &= s, \end{aligned}$$

where the first equation represents the Darcy law, which relates the velocity to the pressure gradient, and the second equation comes from the mass conservation of the fluid. These equations are valid in $\Omega \times T$, where $\Omega \subset R^d$ is a space domain and T is a time interval. The equations should be complemented by proper boundary and initial conditions. The coefficient k represents permeability, c_{pp} storativity depending on fluid and matrix compressibility, s is a source term.

In the case of porous media, which undergo elastic deformations of the matrix, the poroelasticity problem can be described by using the displacement $u = u(x, t)$ as an additionally state variable. The system of equations is then as follows

$$\begin{aligned} -\operatorname{div}(C : \varepsilon(u)) + \alpha \nabla p &= f, \\ K^{-1}v + \nabla p &= 0, \\ \alpha \frac{\partial}{\partial t} \operatorname{div}(u) + \operatorname{div}(v) + c_{pp} \frac{\partial}{\partial t} p &= s. \end{aligned}$$

The first equation is now the Navier- Lamè equation, C is the elasticity tensor, $\sigma_{eff} = C : \varepsilon(u)$ is the effective stress, $\sigma = \sigma_{eff} + \alpha \nabla p$ is the total stress, α is the Biot-Willis coefficient.

The described problems can be formulated variationally in spaces $(H^1(\Omega))^d$ for $u = u(\cdot, t)$, $H(\operatorname{div}, \Omega)$ for $v = v(\cdot, t)$ and $L_2(\Omega)$ for $p = p(\cdot, t)$. The space discretization then used a triple of finite element spaces $U_h \times V_h \times P_h$. The frequently used triple couples continuous piecewise linear elements for u , the Raviart-Thomas elements for v and piecewise constant elements for p . This space discretization provided a differential-algebraic system and its discretization with the backward Euler method leads to the solution of systems in each time step with size τ . Using the row scaling with $(1, \tau, \tau)$, the systems can be written in a scaled symmetric block form, where the blocks correspond to individual variables. They get the following form

$$\mathcal{A}_p = \begin{bmatrix} M & B^T \\ B & -C \end{bmatrix}, \quad \mathcal{A}_{pe} = \begin{bmatrix} A & B_u^T \\ B_u & M & B^T \\ & B & -C \end{bmatrix}.$$

Note that we can assume that the factor τ is included in M and B , i.e. $\tau^{-1}M$ is a positive definite velocity mass matrix, $\tau^{-1}B$ is a matrix representing the divergence of velocity, C is a positive definite pressure mass matrix, which is diagonal for discretization of pressure by piecewise constant finite elements, A is the elasticity stiffness matrix and B_u represents divergence of displacement. The systems with the matrices \mathcal{A}_p and \mathcal{A}_{pe} can be solved iteratively by the Krylov space methods, and due to the symmetry, we can use e.g. the MINRES method.

The aim of this contribution is to introduce and analyse efficient block type preconditioners for systems with matrices \mathcal{A}_p and \mathcal{A}_{pe} , overview the work already done by the authors and outline some new investigations and extensions.

2 Natural Block Diagonal Preconditioners

Let us consider block diagonal preconditioners with the same structure as the matrices \mathcal{A}_p and \mathcal{A}_{pe} . Taking advantage of easy inverse of diagonal C , we can consider the symmetric, positive definite preconditioners (see e.g. [3, 2])

$$\mathcal{P}_p = \begin{bmatrix} M_C & \\ & C \end{bmatrix}, \quad \mathcal{P}_{pe} = \begin{bmatrix} A_C & & \\ & M_C & \\ & & C \end{bmatrix},$$

where $M_C = M + B^T C^{-1} B$ and $A_C = A + B_u^T C^{-1} B_u$ are Schur complements. The algebraic analysis of the preconditioners can use the following theorems.

Theorem 1. *Let*

$$\mathcal{A} = \begin{bmatrix} A_{11} & A_{21}^T \\ A_{21} & -A_{22} \end{bmatrix}, \quad \mathcal{P} = \begin{bmatrix} S & \\ & A_{22} \end{bmatrix},$$

where A_{11} and A_{22} are symmetric positive definite, $S = A_{11} + A_{21}^T A_{22}^{-1} A_{21}$ is Schur complement. Then

$$\sigma(\mathcal{P}^{-1}\mathcal{A}) \subset \left\langle \frac{-1 - \sqrt{5}}{2}, -1 \right\rangle \cup \left\langle \frac{-1 + \sqrt{5}}{2}, 1 \right\rangle.$$

Remark. *The proof of Theorem 1 can be found e.g. in [6]. Theorem 1 can be also generalized for S being only spectrally equivalent to the Schur complement $A_{11} + A_{21}^T A_{22}^{-1} A_{21}$.*

Theorem 2. *Let \mathcal{A} and \mathcal{P} be as above, but S being only spectrally equivalent to the Schur complement,*

$$\xi_0 S \leq A_{11} + A_{21}^T A_{22}^{-1} A_{21} \leq \xi_1 S.$$

Then

$$\sigma(\mathcal{P}^{-1}\mathcal{A}) \subset \left\langle -\frac{1}{2} - \frac{1}{2}\sqrt{1 + 4\xi_1}, -1 \right\rangle \cup \left\langle -\frac{1}{2} + \frac{1}{2}\sqrt{1 + 4\xi_0}, \xi_1 \right\rangle.$$

Theorem 1 can be directly applied to the analysis of the preconditioner \mathcal{P}_p and the preconditioner \mathcal{P}_{pe+} ,

$$\mathcal{P}_{pe+} = \begin{bmatrix} A_C & S_{21}^T & & \\ S_{21} & M_C & & \\ & & & C \end{bmatrix},$$

where $S_{21} = BC^{-1}B_u$.

Theorem 3. *Let $0 \leq \gamma < 1$ be a strengthened Cauchy-Schwarz inequality constant, i.e.*

$$\mathbf{v}^T S_{21} \mathbf{u} \leq \gamma \sqrt{\mathbf{u}^T A_C \mathbf{u}} \sqrt{\mathbf{v}^T M_C \mathbf{v}} \quad \forall \mathbf{u} \in R^{n_u} \equiv u_h \in U_h, \mathbf{v} \in R^{n_v} \equiv v_h \in V_h.$$

Then there is a spectral equivalence

$$(1 - \gamma) \begin{bmatrix} A_C & & \\ & M_C & \end{bmatrix} \leq \begin{bmatrix} A_C & S_{21}^T \\ S_{21} & M_C \end{bmatrix} \leq (1 + \gamma) \begin{bmatrix} A_C & & \\ & M_C & \end{bmatrix}.$$

If c_{el} is a positive constant such that

$$c_{el} \|\operatorname{div}(u_h)\|_{L_2}^2 \leq \langle A \mathbf{u}, \mathbf{u} \rangle \quad \forall \mathbf{u} \in R^{n_u} \equiv u_h \in U_h,$$

then $\gamma^2 \leq (1 + c_{pp}c_{el})^{-1}$. For isotropic elasticity with Lamè constants λ and μ , $c_{el} = \lambda$.

The challenging point of the implementation of the block diagonal preconditioners is the solution of the symmetric positive definite Schur complement systems with matrices M_C and A_C . Note that for piecewise constant finite elements matrices for pressure, these matrices are sparse and can be assembled in the standard element-by-element manner. This enables to use direct as well as many types of preconditioned iterative methods for the solution. For M_C , we suggested to use CG with additive Schwarz preconditioner in [7] and showed that even highly parallelizable one-level Schwarz method can be efficient for parameters corresponding to hardly permeable porous media.

3 Conclusions and Extensions

The considered preconditioners can be also efficiently used for more accurate discretizations. The application for higher order Radau discretization in time is considered in [5], more parallelizable versions of such preconditioning are introduced in [4]. The considered triple of finite elements can suffer from locking and therefore the use of nonconforming discretization for elasticity part was suggested recently, e.g. in [8, 9]. We shall discuss the use of block diagonal preconditioner for this extension. Another extension is for models related to poroelasticity as e.g. Biot-Barenblatt double permeability (see [6]). The preconditioner can be also used for nonlinear models, e.g. for not fully saturated flow described by the Richards equation.

Beside symmetric positive definite block diagonal preconditioners, it is possible to consider block diagonal indefinite and triangular preconditioners in combination with GMRES. Such preconditioners were investigated e.g. in [1, 2].

Acknowledgement: The work was done within the project LQ1602 "IT4Innovations excellence in science" supported by the Ministry of Education, Youth and Sports of the Czech Republic.

References

- [1] O. Axelsson, Unified analysis of preconditioning methods for saddle point matrices. *Numerical Linear Algebra with Applications*, Vol 22 (2015), pp. 233-253
- [2] O. Axelsson, R. Blaheta, P. Byczanski, Stable discretization of poroelasticity problems and efficient preconditioners for arising saddle point type matrices. Volume 15 (2012), Issue 4, pp 191–207
- [3] O. Axelsson, R. Blaheta, P. Byczanski, J. Karátson and B. Ahmad. Preconditioners for regularized saddle point operators with an application for heterogeneous Darcy flow and transport problems. *Journal of Computational and Applied Mathematics* 280 (2015) 141–157
- [4] O. Axelsson, R. Blaheta, T. Lubner: Preconditioners for mixed FEM solution of stationary and nonstationary porous media flow problems. *Large-Scale Scientific Computing, Lecture Notes in Comput. Sci.* 9374, Springer, 2015, pp. 3–14.
- [5] O. Axelsson, R. Blaheta, R. Kohut: Preconditioned methods for high order strongly stable time integration methods with an application for a DAE problem, *Numerical Linear Algebra with Applications*, 2015 (22), pp.930–949
- [6] R. Blaheta, T. Lubner, Algebraic preconditioning for Biot-Barenblatt poroelastic systems. *Spec. issue SNA, Applications of Mathematics*, Vol. 62, No. 6, pp. 561-577, 2017
- [7] R. Blaheta, J. Kruzik, T. Lubner, Schur Complement-Schwarz DD Preconditioners for Non-Stationary Darcy Flow Problems. To appear in *HPCSE 2017, LNCS* Springer
- [8] X. Hu, C. Rodrigo, F.J. Gaspar, L.T. Zikatanov, A nonconforming finite element method for the Biot's consolidation model in poroelasticity. *Journal of Computational and Applied Mathematics* 310 (2017) 143–154
- [9] Q. Hong, J. Kraus, Parameter-robust stability of classical three-field formulation of Biot's consolidation model. Submitted to *ETNA*

A Monolithic Scheme for a Fluid-Poroelastic Structure Interaction Problem

A. Cesmelioglu

The interaction of an incompressible Newtonian fluid with a poroelastic material takes place in important multiphysics problems arising in various applications. For example, blood flow is affected by the porous and deformable nature of the arterial wall and understanding how blood flows in arteries through simulations may be beneficial in biomedical engineering. The model we use is the one in [1, 2]. In this talk, we present a monolithic scheme based on the finite element method and its analysis assuming that the boundary and interface are fixed.

We use the time-dependent incompressible Stokes equations to model the fluid flow in Ω_f :

$$\begin{aligned} \rho_f \mathbf{u}_t - 2\nu_f \nabla \cdot \mathbf{D}(\mathbf{u}) + \nabla p &= \mathbf{f}_f \quad \text{in } \Omega_f \times (0, T), \\ \nabla \cdot \mathbf{u} &= 0 \quad \text{in } \Omega_f \times (0, T), \end{aligned}$$

where \mathbf{u} denotes the velocity vector of the fluid, p denotes the pressure of the fluid, ρ_f denotes the density of the fluid, ν_f denotes the constant fluid viscosity, and \mathbf{f}_f denotes the body force acting on the fluid. The strain rate tensor $\mathbf{D}(\mathbf{u})$ is defined by $\mathbf{D}(\mathbf{u}) = \frac{1}{2} (\nabla \mathbf{u} + (\nabla \mathbf{u})^T)$ and the Cauchy stress tensor is given by $\boldsymbol{\sigma}_f = 2\nu_f \mathbf{D}(\mathbf{u}) - p\mathbf{I}$. To model the poroelastic material in Ω_p we use the Biot equations:

$$\begin{aligned} \rho_s \boldsymbol{\eta}_{tt} - 2\nu_s \nabla \cdot \mathbf{D}(\boldsymbol{\eta}) - \lambda_s \nabla (\nabla \cdot \boldsymbol{\eta}) + \alpha \nabla \phi &= \mathbf{f}_s \quad \text{in } \Omega_p \times (0, T), \\ (s_0 \phi_t + \alpha \nabla \cdot \boldsymbol{\eta}_t) - \nabla \cdot \mathbf{K} \nabla \phi &= f_p \quad \text{in } \Omega_p \times (0, T), \end{aligned}$$

where $\boldsymbol{\eta}$ is the displacement of the structure, ϕ is the pore pressure of the fluid, f_p is the source/sink term and \mathbf{f}_s is the body force. The parameters ν_s and λ_s denote the Lamé constants for the solid skeleton. The density of the saturated medium and the hydraulic conductivity are denoted by ρ_s and \mathbf{K} , respectively. The total stress tensor for the poroelastic structure is given by: $\boldsymbol{\sigma}_p = 2\nu_s \mathbf{D}(\boldsymbol{\eta}) + \lambda_s (\nabla \cdot \boldsymbol{\eta}) \mathbf{I} - \alpha \phi \mathbf{I}$. We assume that \mathbf{K} is a symmetric positive definite tensor such that there exists $K_{\min}, K_{\max} > 0$ satisfying

$$K_{\min} \boldsymbol{\xi} \cdot \boldsymbol{\xi} \leq \boldsymbol{\xi} \cdot \mathbf{K} \boldsymbol{\xi} \leq K_{\max} \boldsymbol{\xi} \cdot \boldsymbol{\xi} \quad \forall \boldsymbol{\xi} \in \bar{\Omega}_p.$$

These equations are coupled via appropriate interface conditions on the interface Γ_I separating Ω_f and Ω_p including continuity for the normal flux, balance of normal stresses, balance of normal components of the stress in the fluid phase and the Beavers-Joseph-Saffman condition which states that the tangential component of the velocity is proportional to the shear stress:

$$\begin{aligned} \mathbf{u} \cdot \mathbf{n}_\Gamma &= (\boldsymbol{\eta}_t - \mathbf{K} \nabla \phi) \cdot \mathbf{n}_\Gamma, \\ \boldsymbol{\sigma}_f \mathbf{n}_\Gamma &= \boldsymbol{\sigma}_p \mathbf{n}_\Gamma, \\ \mathbf{n}_\Gamma \cdot \boldsymbol{\sigma}_f \mathbf{n}_\Gamma &= -\phi, \\ \mathbf{n}_\Gamma \cdot \boldsymbol{\sigma}_f \mathbf{t}_\Gamma &= -\beta (\mathbf{u} - \boldsymbol{\eta}_t) \cdot \mathbf{t}_\Gamma^l, \quad 1 \leq l \leq d-1, d=2,3. \end{aligned}$$

The system is assumed to be at rest initially and simple boundary conditions are chosen. This problem is a coupled system of mixed hyperbolic-parabolic type and inherits all the difficulties mathematically and numerically involved in the standard fluid-structure interaction and fluid-porous media flow coupling problems.

We first reduce this second order in time problem to first order by introducing an additional variable $\zeta = \eta_t$. We present a weak formulation whose analysis can be done as in [3] and based on this weak formulation, we derive a monolithic scheme using the backward Euler method to discretize in time and the finite element method to discretize in space. For the discretization in Ω_f , we use finite element spaces that satisfy the inf-sup condition. We denote by $\mathbf{X}_f^h, Q_f^h, \mathbf{X}_p^h, \mathbf{X}_p^h$ and Q_p^h the finite element spaces to approximate $\mathbf{u}, p, \eta, \zeta, \phi$, respectively.

Let $N \in \mathbb{N}$ and $\Delta t = T/N$ be the time step size. Define $t_i = i\Delta t, i = 0, \dots, N$ and $\mathcal{D}_{\Delta t} g^n := \frac{g^n - g^{n-1}}{\Delta t}$. Then the scheme is as follows: First we find $(\mathbf{u}_h^0, \eta_h^0, \zeta_h^0, \phi_h^0) \in \mathbf{X}_f^h \times \mathbf{X}_p^h \times \mathbf{X}_p^h \times Q_p^h$ by interpolating the initial conditions. Then for all $1 \leq n \leq N$, we seek $(\mathbf{u}_h^n, p_h^n, \eta_h^n, \zeta_h^n, \phi_h^n) \in \mathbf{X}_f^h \times Q_f^h \times \mathbf{X}_p^h \times \mathbf{X}_p^h \times Q_p^h$, such that

$$\begin{aligned} & \rho_f(\mathcal{D}_{\Delta t} \mathbf{u}_h^n, \mathbf{v})_{\Omega_f} + a_f(\mathbf{u}_h^n, \mathbf{v}) + b_f(\mathbf{v}, p_h^n) + \rho_s(\zeta_h^n - \mathcal{D}_{\Delta t} \eta_h^n, \mathbf{tau})_{\Omega_p} + \rho_s(\mathcal{D}_{\Delta t} \zeta_h^n, \xi)_{\Omega_p} \\ & + a_e(\eta_h^n, \xi) - b_e(\xi, \phi_h^n) + s_0(\mathcal{D}_{\Delta t} \phi_h^n, r)_{\Omega_p} + b_e(\mathcal{D}_{\Delta t} \eta_h^n, r) + a_d(\phi_h^n, r) + \langle \phi_h^n \mathbf{n}_\Gamma, \mathbf{v} - \xi \rangle_{\Gamma_I} \\ & + \sum_{l=1}^{d-1} \langle \beta(\mathbf{u}_h^n - \mathcal{D}_{\Delta t} \eta_h^n) \cdot \mathbf{t}_\Gamma^l, (\mathbf{v} - \xi) \cdot \mathbf{t}_\Gamma^l \rangle_{\Gamma_I} + \langle (\mathcal{D}_{\Delta t} \eta_h^n - \mathbf{u}_h^n) \cdot \mathbf{n}_\Gamma, r \rangle_{\Gamma_I} \\ & = - \langle P_{in}^n \mathbf{n}_f, \mathbf{v} \rangle_{\Gamma_f^n} + (\mathbf{f}_f^n, \mathbf{v})_{\Omega_f} + (\mathbf{f}_s^n, \xi)_{\Omega_p} + (f_p^n, r)_{\Omega_p}, \\ b_f(\mathbf{u}_h^n, q) & = 0, \quad \forall (\mathbf{v}, q, \xi, \chi, r) \in \mathbf{X}_f^h \times Q_f^h \times \mathbf{X}_p^h \times \mathbf{X}_p^h \times Q_p^h \end{aligned}$$

where

$$\begin{aligned} a_f(\mathbf{v}, \mathbf{w}) & = 2\nu_f(\mathbf{D}(\mathbf{u}), \mathbf{D}(\mathbf{w}))_{\Omega_f}, \forall \mathbf{v}, \mathbf{w} \in \mathbf{X}_f^h, \\ b_f(\mathbf{v}, q_f) & = -(q_f, \nabla \cdot \mathbf{v})_{\Omega_f}, \forall \mathbf{v} \in \mathbf{X}_f^h, \forall q_f \in Q_f^h, \\ a_e(\eta, \xi) & = (2\nu_s \mathbf{D}(\eta), \mathbf{D}(\xi))_{\Omega_p} + (\lambda_s \nabla \cdot \eta, \nabla \cdot \xi)_{\Omega_p}, \forall \eta, \xi \in \mathbf{X}_p^h, \\ b_e(\xi, q_p) & = \alpha(q_p, \nabla \cdot \xi)_{\Omega_p}, \forall \xi \in \mathbf{X}_p^h, q_p \in Q_p^h, \\ a_d(q_p, \psi) & = (\mathbf{K} \nabla q_p, \nabla \psi)_{\Omega_p}, \forall q_p, \psi \in Q_p^h. \end{aligned}$$

We prove that there exists a unique discrete solution $\{(\mathbf{u}_h^n, p_h^n, \eta_h^n, \zeta_h^n, \phi_h^n)\}_{n \geq 0}$ and the method is stable, that is, the discrete solution is bounded by the data of the problem $\mathbf{f}_f, \mathbf{f}_s, f_p, \rho_f, \nu_f, \rho_s, \nu_s, \alpha, s_0, \mathbf{K}$. The error is proved to be optimal in the sense that if we use polynomials of degree k_1 for \mathbf{u}_h , $k_1 - 1$ for p_h , k_2 for η_h, ζ_h and $k_2 - 1$ for ϕ_h , then the convergence is of order $\mathcal{O}(h^{2k_1} + h^{2k_2} + (\Delta t)^2)$ where h is the mesh size.

References

- [1] M. Bukač, I. Yotov, R. Zakerzadeh, P. Zunino, *Partitioning strategies for the interaction of a fluid with a poroelastic material based on a Nitsche's coupling approach*, Comput. Methods Appl. Mech. Engrg. 292 (2015) 138–170, Special

Issue on Advances in Simulations of Subsurface Flow and Transport (Honoring Professor Mary F. Wheeler).

- [2] A. Cismelioglu, H. Lee, A. Quaini, K. Wang, S.-Y. Yi, *Optimization-based decoupling algorithms for a fluid-poroelastic system*, in: S.C. Brenner (Ed.), *Topics in Numerical Partial Differential Equations and Scientific Computing*, Springer New York, New York, NY, 2016, pp. 79–98.
- [3] A. Cismelioglu, *Analysis of the coupled Navier-Stokes/Biot problem*, *J. Math. Anal. Appl.* 456 (2017) 970–991.

Computational Nano-physics – Monte Carlo Approach

I. Dimov

The Wigner equation is a full quantum model capable of capturing the relevant physics needed for the simulation of nano-devices. There have been a number of recent publications that deal with the numerical approximation of observables related to the Wigner equation using probabilistic techniques and most notably Monte-Carlo methods based on branching particle systems. In this work we consider and study in detail the Signed Particle Wigner Monte-Carlo method, which is devised for the numerical approximation of observables related to the Wigner equation.

Convergence of class of Monte Carlo methods dealing with observables in Quantum Physics is analyzed. We deal with the numerical approximation of observables related to the Wigner equation using probabilistic techniques based on branching particle systems. We answer several questions about the behavior of the algorithm and demonstrate theoretically why almost always is not stable and how to deal with this instability. Our work relies exclusively on probabilistic techniques and the estimates related to the proposed algorithms can be seen as sharpening of the more general study of stochastic algorithms for the Wigner equation. The work also summarizes the formulation of the Wigner equation as an operator equation in suitable L_2 spaces.

Numerical Approximation of Fractional Spectral Elliptic Operators

B. Duan, R.D. Lazarov, J.H. Pasciak

We shall discuss methods and algorithms for approximately solving the linear algebraic systems $L_h^\alpha u_h = v_h$, $0 < \alpha < 1$, for $u_h, v_h \in H_h$, H_h a finite dimensional Hilbert space. Such problems arise in finite element or finite difference approximations of problems $L^\alpha u = v$ with fractional powers of second order elliptic operators L . The algorithms are based on the method of Vabishchevich, that related the algebraic problem to a solution of a time-dependent parabolic type equation on the interval $[0, 1]$.

We develop and study two algorithms based on diagonal Padé approximation of the corresponding solution operator. The first one uses geometrically graded meshes in order to compensate for the singular behavior of the solution for t close to 0 for non-smooth data v . The second algorithm uses uniform in t meshes, but requires smoothness of the data v in order to retain optimal convergence rate. For both methods we estimate the error in terms of the number of time steps and the regularity of the data. Finally, we report some numerical experiments of finite element approximation of second order elliptic problems in one and two spatial dimensions.

The work of R. Lazarov was supported in parts by the Bulgarian NSF Grant DN 12/1 and USA NSF-DMS Grant #1620318

Deflation Methods Made Possible

S.-E. Ekström, M. Neytcheva

It is well known that the convergence of Krylov subspace iterative solution methods is severely hampered by the presence of small eigenvalues (e.g., [1]). Techniques how to avoid or diminish this effect are also well-known, under the names *bordering*, *augmenting* or *deflation*. All those techniques are based on the assumption that we know the smallest eigenvalues and their corresponding eigenvectors, or, at least some good enough approximations of those. At the same time, we are aware that, in general, to compute or estimate a number of eigenvalues and eigenvectors is not an easy task, in particular, for large size matrices. Note that not only the eigenpair computation is time consuming. It might be practically infeasible when we are aiming to determine a number of interior eigenvalues or all eigenvalues in an interval.

Even though the above difficulties are widely acknowledged, deflation methods have been recently advertized in the context of High Performance Computing (HPC) and in particular, in the search of iterative solution methods, suitable for exascale computer platforms.

In this talk we show that for certain classes of problems we are able to compute the eigenvalues exactly or within machine precision, without even constructing the matrices explicitly. The target matrices arise from discretizations of partial differential models, discretized using local methods, such as Finite Elements (FEM), Finite Differences (FD), Finite Volumes and Iso-geometric Analysis (IgA). For some problems it is possible to construct also explicitly the eigenvectors. At this stage of development of the techniques the discretization meshes are regular and tensor-based. The latter, although seen as a disadvantage for some problems, is also attractive for HPC applications.

To put the presentation in some context, consider the solution of a linear system of equations

$$A\mathbf{x} = \mathbf{b},$$

where $A \in R^{n \times n}$ is a large and sparse nonsingular matrix. We also assume that A is symmetric and positive definite. Let $\{\lambda_i\}$ and $\mathbf{v}^{(i)}$, $i = 1, 2, \dots, n$ be the eigenvalues of A and their corresponding eigenvectors. We start with the assumption that we know the smallest eigenvalues of A , respectively, their eigenvectors,

Recall one of the first frameworks that eliminate the influence of several nearly zero eigenvalues on the condition number of A , say p of them. Consider the so-called *bordering* method, described in [2]. We construct a larger matrix by bordering the original one with the corresponding number of columns and rows. The eigenvalues are then perturbed as stated in the next theorem.

Theorem Let A be of order $n \times n$ and V_p of order $n \times p$ where $p < n$. Consider the augmented system

$$\tilde{A} = \begin{bmatrix} A & -AV_p \\ -V_p^T A & V_p^T AV_p \end{bmatrix}.$$

Then

- (a) \tilde{A} has p zero eigenvalues. The remaining eigenvalues $\tilde{\lambda}_i$ are equal to those of $(I + V_p V_p^T)A$.
- (b) If A is s.p.d., then to every eigenvalue λ_i of A there exists $\tilde{\lambda}_i$ such that $\tilde{\lambda}_i \geq \lambda_i$.
- (c) If A is nonsingular and symmetric and $V_p = [\alpha_1 \mathbf{v}^{(1)}, \dots, \alpha_p \mathbf{v}^{(p)}]$, where $\mathbf{v}^i, i = 1, 2, \dots, p$ are normalized eigenvectors of A , then the nonzero eigenvalues of \tilde{A} equal $\tilde{\lambda}_i = (1 + \alpha_i^2)\lambda_i, i = 1, 2, \dots, p$ and $\tilde{\lambda}_i = \lambda_i, i = p + 1, \dots, n$.
- (d) the smallest effective condition number of \tilde{A} , i.e. $\lambda_{max}((I + V_p^T V_p)A) / \lambda_{min}((I + V_p^T V_p)A)$ obtained by bordering with p vectors is $\lambda_n / \lambda_{p+1}$.

So, we eliminate the p smallest eigenvalues and make the matrix exactly singular. At the same time we do not perturb the rest of the spectrum significantly, ensuring a much smaller effective condition number of \tilde{A} . The resulting singular system is expected to be efficiently solved by a Krylov subspace method.

Of course, in practice we would not advocate to solve a singular matrix of larger size than A . The above framework is shown to be equivalent to solve $(I + V_p^T V_p)A$, which is of the same size as A . See some earlier studies in [3].

The same idea to *deflate* the small eigenvalues has been incorporated in particular implementations of some of the most used Krylov subspace iteration methods, such as the Conjugate Gradient method, (deflated CG), BiCG, GMRES, e.g. [4]. Deflation techniques are advocated also in the context of communication-avoiding techniques and pipelined versions of various iterative solvers.

In all deflation-related studies the question how to compute or approximate the eigenpairs in question, remains a major issue. There are various results on constructing approximate subspaces to be used instead of the exact eigenvectors, which we do not consider here.

Since many years, a (seemingly unrelated to the deflation techniques) scientific theory has been developed, namely, the so-called 'Generalized Locally Toeplitz' (GLT) sequences, cf. [6]. For certain classes of structured matrices, much richer than the classical Toeplitz matrices, GLT offers the possibility to associate an analytical function to (a sequence of) matrices, referred to as the symbol of the matrices. Sampling the symbol gives an information about the spectrum of the corresponding matrix, namely, a curve (for s.p.d. matrices) on which all eigenvalues are located, except for possibly a finite number of outliers. Until recently, it was not known however, where exactly on that curve the exact eigenvalues are located. A significant improvement in this direction is due to the work of the first author and coauthors, cf. [7], providing a methodology to compute exactly (or up to machine accuracy) *all* eigenvalues of a matrix of the considered class, only based on the symbol, in a cheap and easy to implement way.

The knowledge of the exact eigenvalues opens the door to reviving various methods, based on eigenvalue information, in particular, deflation techniques. The question

regarding the eigenvectors is more difficult, however, some approaches are applicable, again, claimed to be profitable in HPC computations.

We show the idea on how to compute the eigenvalues based on the matrix symbol and illustrate the effect on solving the linear systems with deflated methods on problems, discretized by FD, FEM and IgA.

References

- [1] O. Axelsson, *Iterative Solution Methods*, Cambridge University Press, 1994.
- [2] D.K. Faddeev, V.N. Fadееva, *Computational Methods of Linear Algebra*, Freeman, San Francisco, 1963.
- [3] O. Axelsson, M. Neytcheva, B. Polman, The bordering method as a preconditioning method, *Vestnik Moskovskogo Universiteta, Seria 15, V'chisl. Math. Cybern.*, 1 (1996), 3–25.
- [4] Y. Saad, M. Yeung, J. Erhel, and F. Guyomarc'h, A Deflated Version of the Conjugate Gradient Algorithm, *SIAM J. Sci. Comput.*, 21(2006), 1909–1926.
- [5] E. Carson, N. Knight, J. Demmel, An efficient deflation technique for the communication-avoiding conjugate gradient method, *ETNA*, 43 (2014), 125–141.
- [6] C. Garoni, S. Serra-Capizzano, *Generalized Locally Toeplitz Sequences: Theory and Applications*, Springer, 2017.
- [7] S.-E. Ekström, *Matrix-Less Methods for Computing Eigenvalues of Large Structured Matrices*, Ph.D. Thesis, Uppsala University, 2018, https://www.researchgate.net/publication/324653041_Matrix-Less_Methods_for_Computing_Eigenvalues_of_Large_Structured_Matrices.

Reliable Numerical Models and Their Applications

I. Faragó, R. Horváth, M. Mincsovics, R. Mosleh, F. Dorner

1 Introduction and Motivation

Mathematical models are efficient tools of modelling of different phenomena. In the modelling process we formulate these phenomena in the language of mathematics. Typically, the construction of the models is realized with the modelling chain

physical / biological model \rightarrow continuous model \rightarrow discrete (numerical) model.

In order to have an adequate model, it is almost obvious that the continuous model and the numerical model on some fixed mesh should preserve the basic (scientifically motivated) qualitative properties of the original phenomenon. Such models are called qualitatively adequate, or, in short reliable models.

In the following we examine the preservation of the qualitative properties for the numerical models. In some earlier works, we have investigated the discrete model of the heat conduction problem (e.g. [3, 4]). For the heat conduction process the main and physically motivated characteristic properties are the non-negativity preservation, the maximum / minimum principle and the contractivity in the maximum norm. Obviously, these requirements are quite natural and they are motivated by the basic physical principles. Therefore the discrete models for the heat conduction process should have the discrete analogue of these properties. One can show (e.g. [3, 4]) that the connection between the above discrete qualitative properties in the discrete one-step models is the following

maximum principle \Leftrightarrow non-negativity preservation \Rightarrow contractivity

In these papers we formulated the conditions under which the discrete models are reliable. Typically these conditions result in some restriction for the choice of the discretization parameters, namely, for fixed space discretization there are bounds for the time-discretization step-size.

In the following we focus our attention on some discrete mathematical models of the biology, namely we consider some discrete epidemic models and we will investigate their qualitative properties. This problem is related to the following phenomenon. The modeling of infectious diseases is a tool which has been used to study the mechanisms by which diseases spread, to predict the future course of an outbreak and to evaluate strategies to control an epidemic.

Our basic aim is to investigate the corresponding discrete models from qualitative point of view. As the qualitative property, for the first problem we investigate the mass preservation, monotonicity and the non-negativity preservation properties, while, for the second problem we analyse the energy-preservation property.

2 Qualitative Properties of Discrete Epidemic Models

Our natural demand is to predict the spatial motion of diseases, and to prevent or to curb epidemics. Mathematical models are very effective tools of the investigation of disease propagations. This is why a number of mathematical models have been constructed and investigated in the literature, see e.g. [1]. The investigations started almost one hundred years ago with the system of ordinary differential equations model given in [7]. This model is a so-called compartmental model, where the population is divided into disjoint groups according to the members' relation to the disease, and the time-dependent function of the number of the members in each group is determined with the solution of the system of the ordinary differential equations. The most typical compartments are as follows: susceptibles (members that can be infected), infectives (members that can pass on the disease to others) and recovered (members that have recovered from the disease). Compartmental models describe only the number of the members in each compartment as a function of time, thus they are not able to give spatial information about the disease. Models in the form of systems of partial differential equations must be constructed to incorporate spatial dependence into the models, e.g. [1]. In the case of these models mainly the stability of the stationary states and the pattern formation is investigated and little attention is paid to the qualitative properties of the numerical solutions. When the birth and the death of the members are not taken into the account then the main qualitative properties of the disease propagation process are the mass conservation, non-negativity preservation and monotonicity of the number of the susceptibles and the recovered members.

Because of the spatial dependence, the above properties will be formulated for the discrete density functions of the compartments. The densities of the members in the susceptible, infective and recovered compartment are denoted by $S(x, t)$, $I(x, t)$ and $R(x, t)$, respectively. This means, for example, that when we integrate the function $I(x, t)$ on a spatial domain Ω then we get the number of infective members in the domain Ω at the time instant t .

Now we construct discrete spatial disease propagation models. We will give the conditions of the validity of the above qualitative properties for these models. We define the problem on the cubical domain $[0, L]^D$ ($L > 0$). Here D is the spatial dimension of the problem ($D = 1$ and 2 are the important dimensions from the practical point of view). In the discrete model we apply a uniform spatial grid

$$\omega_h = \{(x_1, \dots, x_D) \in [0, L]^D \mid x_k \in \{0, h, 2h, \dots, Mh, (M+1)h\}, k = 1, \dots, D, \\ h = L/(M+1), M \in \mathbb{M}^+\}$$

and a positive time step $\tau > 0$. The sub-population density functions S , I and R are approximated, respectively, by the grid functions S^n , I^n and R^n at the n th time level $t = n\tau$. The values of the grid functions with $n = 0$ are known from the initial conditions, moreover the values of the functions are equal zero in the boundary points. In order to calculate the values of the grid functions in the inner points, we reshape the

values of the grid-functions into column vectors using the usual column-wise indexing. In this way, we obtain the column vectors s^n, i^n and $r^n \in \mathbb{R}^{M^D}$. We consider two discrete models. Both are the combinations of implicit and explicit time-stepping methods. The first discrete model

$$\begin{aligned}\frac{s^{n+1} - s^n}{\tau} &= -s^{n+1} \circ p^n \\ \frac{i^{n+1} - i^n}{\tau} &= s^{n+1} \circ p^n - bi^{n+1}, \\ \frac{r^{n+1} - r^n}{\tau} &= bi^{n+1},\end{aligned}\tag{1}$$

is an implicit-explicit (IMEX) discretization of the continuous disease propagation model

$$\begin{aligned}S'_t &= -S(\vartheta I + \varphi \Delta_D I), \\ I'_t &= S(\vartheta I + \varphi \Delta_D I) - bI, \\ R'_t &= bI,\end{aligned}\tag{2}$$

where $p^n = \vartheta i^n + (\varphi/h^2)Q_D i^n$ and Q_D/h^2 is the discretization matrix of the D dimensional Laplace operator Δ_D . We suppose homogeneous Dirichlet boundary conditions. The parameters b (the recovery parameter), ϑ and φ are positive numbers. The last two values are calculated from a weighting function that gives the locality of the infection (e.g. [2, 5]). The matrix Q_D can be formed as follows. Let us define the tridiagonal matrix $Q = \text{tridiag}(1, -2, 1) \in \mathbb{R}^{M \times M}$. Then, if $D = 1$ then $Q_D = Q$ and in the case of $D = 2$ we have $Q_D = I_M \otimes Q + Q \otimes I_M$, where $I_M \in \mathbb{R}^{M \times M}$ is the identity matrix and \otimes denotes the Kronecker product.

The conditions of the qualitative properties of the scheme (1) can be given as follows.

Theorem 4. *Assume that at the initial state $s^0 \geq 0$, $i^0 \geq 0$, $r^0 \geq 0$, and $p^0 \geq 0$, moreover suppose that*

$$\tau \leq \begin{cases} 1/((2D\varphi/h^2 - \vartheta)M_{\max}), & \text{if } h < h^*, \\ \text{arbitrary}, & \text{if } h \geq h^*, \end{cases}\tag{3}$$

where $h^* = (2D\varphi/\vartheta)^{1/2}$ and $M_{\max} = \max(s^0 + i^0 + r^0)$. Then the scheme (1) satisfies the mass conservation property, the non-negativity preservation property and the monotonicity property. The first and the third property is meant pointwise.

The second discrete model has the form

$$\begin{aligned}\frac{s^{n+1} - s^n}{\tau} &= \frac{d_S}{h^2} \bar{Q}_D s^{n+1} - ki^n \circ s^n, \\ \frac{i^{n+1} - i^n}{\tau} &= \frac{d_I}{h^2} \bar{Q}_D i^{n+1} + ki^n \circ s^n - bi^n, \\ \frac{r^{n+1} - r^n}{\tau} &= \frac{d_R}{h^2} \bar{Q}_D r^{n+1} + bi^n.\end{aligned}\tag{4}$$

This model comes from an IMEX discretization of the continuous disease propagation model [1]

$$\begin{aligned} S'_t(x, t) &= d_S \Delta_D S(x, t) - kI(x, t)S(x, t), \\ I'_t(x, t) &= d_I \Delta_D I(x, t) + kI(x, t)S(x, t) - bI(x, t), \\ R'_t(x, t) &= d_R \Delta_D R(x, t) + bI(x, t). \end{aligned} \quad (5)$$

Here d_S , d_I and d_R are the diffusion parameters of the sub-populations, moreover k and b are positive constants. In the discrete model \bar{Q}_D/h^2 denotes the second order approximation of the Laplace operator taking into the account the homogeneous Neumann boundary condition. The matrix \bar{Q}_D can be constructed by using the tridiagonal matrix

$$\bar{Q} = \begin{bmatrix} -2 & 2 & 0 & \cdots & \\ 1 & -2 & 1 & 0 & \cdots \\ & & \cdots & & \\ \cdots & 0 & 1 & -2 & 1 \\ & \cdots & 0 & 2 & -2 \end{bmatrix} \in \mathbb{R}^{(M+2) \times (M+2)}. \quad (6)$$

Then, if $D = 1$ then $\bar{Q}_D = \bar{Q}$ and in the case of $D = 2$ we have $\bar{Q}_D = I_{M+2} \otimes \bar{Q} + \bar{Q} \otimes I_{M+2}$ [8]. The qualitative properties of the model can be guaranteed by the following theorem.

Theorem 5. *Let us consider the IMEX scheme for the system (5) with homogeneous Neumann boundary condition. The mass conservation property is satisfied without any condition, moreover the condition*

$$\tau \leq \min \left\{ \frac{h^D}{2kDN^*}, \frac{1}{b} \right\} \quad (7)$$

implies the non-negativity property and the monotonicity property. N^ is the approximation of the integral of $S + I + R$ computed with the trapezoidal rule at the initial time instant. The mass conservation and the monotonicity properties are related to the number of the members in the whole domain.*

The next model is constructed to simulate the spreading of malaria. As it is known, malaria is an infectious disease caused by the Plasmodium parasite and transmitted between humans through bites of female Anopheles mosquitoes. A mathematical model describes the dynamics of malaria and human population compartments in terms of mathematical equations and these equations represent the relations between relevant properties of the compartments. The discrete model of the malaria, which is investigated, is the following

$$\begin{aligned} x_{n+1} &= x_n + \tau (\alpha y_n (1 - x_n) - r x_n) \\ y_{n+1} &= y_n + \tau (\beta x_n (1 - y_n) - \mu y_n), \end{aligned} \quad (8)$$

which is the explicit discretization by Euler method of the continuous propagation model of malaria, first given by Ross

$$\begin{aligned}\dot{x}(t) &= \alpha y(t)(1 - x(t)) - rx(t) \\ \dot{y}(t) &= \beta x(t)(1 - y(t)) - \mu y(t)\end{aligned}\tag{9}$$

using the step-size τ . In these models the parameters α, β, r and μ are given biological parameters and the unknown functions $x(t)$ and $y(t)$ (and their discretization) yield the density of infected human and mosquitoes, respectively. Therefore, the natural requirement for the models is that they are non-negative and not greater than one. This qualitative properties of the model can be guaranteed by the following theorem.

Theorem 6. *Let us assume that in the discrete model (8) the time discretization step-size satisfies the condition*

$$\tau \leq \min \left\{ \frac{1}{\alpha}, \frac{1}{\beta}, \frac{1}{\mu}, \frac{1}{r} \right\}.\tag{10}$$

Then the dynamical system (8) is invariant with respect to the set $[0, 1]$, i.e., if $x_0, y_0 \in [0, 1]$ then $x_n, y_n \in [0, 1]$ for any n , too.

We can construct another discrete model for the Ross-system (9), by using implicit Euler method. For such approach we build the following model

$$\begin{aligned}x_{n+1} &= x_n + \tau (\alpha y_{n+1}(1 - x_{n+1}) - rx_{n+1}) \\ y_{n+1} &= y_n + \tau (\beta x_{n+1}(1 - y_{n+1}) - \mu y_{n+1}).\end{aligned}\tag{11}$$

Then the following statement holds.

Theorem 7. *The discrete model (11) is unconditionally invariant with respect to the set $[0, 1]$, i.e., if $x_0, y_0 \in [0, 1]$ then $x_n, y_n \in [0, 1]$ for any n , too.*

References

- [1] Capasso, V., Mathematical Structures of Epidemic Systems, Lecture notes in biomathematics 97, Springer (2008).
- [2] Faragó I., Horváth R., On a spatial epidemic propagation model. Progress in Industrial Mathematics at ECMI 2014 (Taormina, Italy, June 2014), Editors: Russo, G., Capasso, V., Nicosia, G., Romano, V., Springer, ISBN 978-3-319-23412-0, 2017.
- [3] Faragó I., Horváth R., Discrete maximum principle and adequate discretizations of linear parabolic problems. SIAM Scientific Computing. 28, 2313–2336. (2006).
- [4] Faragó I., Horváth R., Continuous and discrete parabolic operators and their qualitative properties. IMA Numerical Analysis. 29, 606–631. (2009).
- [5] Faragó I., Horváth R., On some qualitatively adequate discrete space-time models of epidemic propagation. J. Comput. Appl. Math. 293, 45–54 (2016).

- [6] Faragó I., Horváth R., Qualitative Properties of Numerical Solutions of some PDE Models of Disease Propagation. *J. Comput. Appl. Math.*, to appear.
- [7] Kermack, W. O., McKendrick, A. G.: A contribution to the mathematical theory of epidemics. In: *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* 115 (772) pp. 235–240 (1927).
- [8] Thomas, J.W., *Numerical Partial Differential Equations, Finite Difference Methods*, Texts in Applied Mathematics, 22, Springer (1995).

Vertical Structure of Atmospheric Composition Fields over Bulgaria

G. Gadzhev, K. Ganev

Abstract

The numerical simulations for the vertical structure of atmospheric composition fields over Bulgaria have been performed using the US EPA Model-3 system as a modelling tool for 3D simulations and the system nesting capabilities were applied for downscaling the simulations to a 9 km resolution over Bulgaria. The national emission inventory was used as an emission input for Bulgaria, while outside the country the emissions are from the TNO high resolution inventory. The air pollution pattern is formed as a result of interaction of different processes, so if know the contribution of each, for different meteorological conditions and given emission spatial configuration and temporal behavior is very helpful for understanding the atmospheric composition and pollutants behavior. Numerically obtained characteristics for the vertical structure of the atmospheric composition will be demonstrated in the paper.

Introduction

The parameters of the atmosphere have key impact on quality of life and human health. Because of this, quite naturally, the surface air quality is mostly studied. From the other hand the atmospheric composition fields are formed as a result of complex interaction of processes with different temporal and spatial scales — from global to synoptic to a chain of local scales. A very significant role in the formation of air pollution pattern play also the atmospheric turbulence and the atmospheric boundary layer processes. The impact of complex terrain phenomena is also important. The picture becomes even more complex in urban environment, where human activities leads to the formation of specific urban climate, urban heat island and urban boundary layer with complex structure.

The incredible diversity of dynamic processes, the complex chemical transformations of the compounds and complex emission configuration together lead to the formation of a complex vertical structure of the atmospheric composition. The detailed analysis of this vertical structure with it temporal/spatial variability jointly with the atmospheric dynamics characteristics can enrich significantly the knowledge about the processes and mechanisms, which form air pollution, including near earth surface. The present paper present first results of a study, which aims at performing reliable, comprehensive and detailed analysis of the atmospheric composition fields 3D structure and its connection with the processes, which lead to their formation.

1 Methodology

The study is based on ensemble of computer simulations of the atmospheric composition fields in Bulgaria. The simulations have been performed using US EPA Model-3 system as modeling tools for 3D simulations: Meteorological model WRF [5], Atmosphere Composition Model CMAQ [1] and Emission model SMOKE [2]. The NCEP Global Analysis Data meteorological background with $1^\circ \times 1^\circ$ resolution was used. The models nesting capabilities were applied to downscale the simulations to 9 km for Bulgaria. The TNO high resolution emission inventory [7] and National emission inventory as emission input for Bulgaria have been used. More detailed description of the experiments can be seen in [3],[4].

The computer resource requirements of the modeling system are very big [6] and the numerical experiments were organized in effective HPC environment. The calculations were implemented on the Supercomputer System “Avitohol” at IICT-BAS.

The simulations, performed day by day for 7 years (2008–2014), produced ensemble, comprehensive enough as to provide statistically reliable assessment of the atmospheric composition climate. By averaging over the ensemble the “typical” seasonal and annual pollution concentration fields were constructed with their spatial variability and diurnal course. Some characteristics of the concentration fields vertical distribution – center of masses, vertical dispersion, skewness and kurtosis, vertical mean and maximal concentration and the maximal concentration level have been calculated for the so constructed “typical” seasonal and annual pollution concentration fields.

2 Results

The center of masses horizontal distribution for the annually averaged Ozone (O_3) and Nitrogen Dioxide (NO_2) at 12:00 GMT are shown on Figure 1. What can immediately be seen from the figure is that the center of masses is significantly horizontally heterogeneous, which reflects the complex interaction of different dynamic and chemistry processes as well as the emission source configuration. It should be also be mentioned that the center of masses for O_3 is much higher than the one for NO_2 . This is a reflection of the already known [3],[4] fact that the surface O_3 in Bulgaria is to great extend due to transport from abroad and/or from above.

Figure 1: Centre of masses horizontal distribution over Bulgaria for the annually averaged NO_2 and O_3 at 12:00 GMT

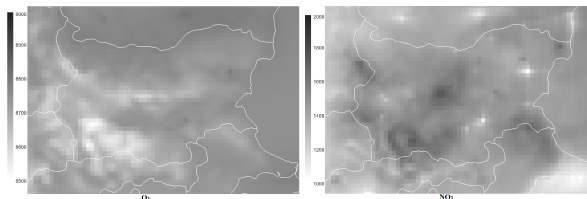


Figure 2: Diurnal course of the annually averaged profiles of O_3 and NO_2 (a), together with the surface, maximal and vertically mean concentrations [$\mu\text{g}/\text{m}^3/\text{h}$] (b), centers of masses, vertical dispersions and maximal concentrations levels [m] (c).

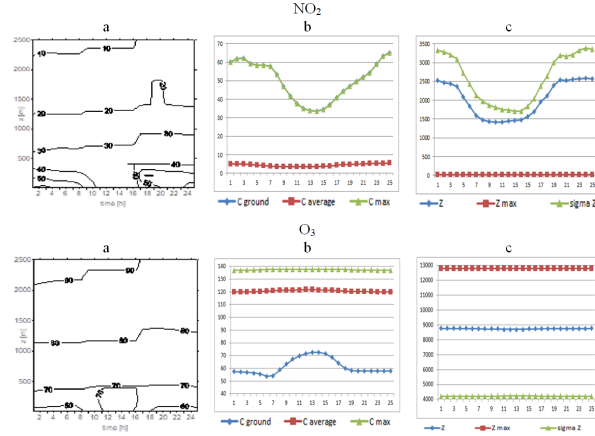


Figure 2 demonstrates the diurnal course of the annually averaged profiles of O_3 and NO_2 , together with the surface, maximal and vertically mean concentrations, the centres of masses, vertical dispersions and maximal concentrations levels. The main vertical/diurnal pollutant behavior features that can be seen from the figure are the following: (i) the diurnal course of the characteristics is generally well manifested; (ii) the NO_2 maximal concentrations always appear in the ground layer, while for the O_3 they are at much higher levels — yet another manifestation of the above mentioned O_3 origin in Bulgaria; (iii) as should be expected the surface NO_2 has local minimum, while the O_3 has a local maximum around and after noon (the joint effect of intensified turbulent transport and photochemistry); (iv) as it should be expected the centre of masses and the vertical dispersion for the NO_2 have similar temporal behavior and close values; (v) against the expectations the centre of masses and the vertical dispersion have local minimum during the warmer part of the day. This is a demonstration that in the environment of a complex 3D dynamics, chemical transformations and emission sources at different levels the admixture propagation can be different from the case of a classic instantaneous point source in conventional boundary layer; (vi) there is practically no diurnal course of all the O_3 characteristics, except the surface concentration, which is another evidence that the O_3 in the major part of the column is not a subject of local processes; the vertical mean concentration of the NO_2 is much smaller than the surface one, while for the O_3 it is higher than the surface one and closer to the maximal. This can be followed also in the vertical distribution skewness (not demonstrated in the paper), which is positive for all the compounds except O_3 .

3 Conclusion

Due to volume limitations only few of the calculated pollution vertical distribution characteristics are shown in the present paper. Nevertheless they are a good demonstration of the complexity of the atmospheric composition 3D structure and of the joint impact of processes of different nature and spatial/temporal scales. The horizontal resolution of the performed computer simulations is rather coarse (9 km), so the simulations do not reflect the role of many local scale phenomena. Performing and analysis of simulations with finer resolution (down to urban scales) will be the task for future work.

Acknowledgment

The present work is supported by: the Bulgarian National Science Fund (grant DCVP-04/2/13.12.2016), EU -H2020 project 675121(project VI-SEEM), Program for career development of young scientists, BAS. Deep gratitude to US EPA and US NCEP, EMEP and TNO for providing free-of-charge data and software.

References

- [1] CMAQ user guide, (2006). [Online] Available from: <https://www.cmascenter.org/help/documentation.cfm?model=cmaq&version=4.6>
- [2] CEP (2003) Sparse Matrix Operator Kernel Emission (SMOKE) Modeling System, University of Carolina, Research Triangle Park, North Carolina.
- [3] Gadzhhev, G., Ganev, K., Miloshev, N., Syrakov, D., Prodanova, M., (2013) Numerical Study of the Atmospheric Composition in Bulgaria, *Computers and Mathematics with Applications* 65, p. 402-422.
- [4] Gadzhhev, G., Ganev, K., Prodanova, M., Syrakov, D., Atanasov, E., Miloshev, N. (2013) Multi-scale Atmospheric Composition Modelling for Bulgaria. *NATO Science for Peace and Security Series C: Envir. Security*, 137, pp. 381-385.
- [5] Shamarock W., Klemp J., Dudhia J., Gill D., Barker D., Wang W., Powers J., (2007) A description of the Advanced Research WRF Version 2, http://www.mmm.ucar.edu/wrf/users/docs/arw_v2.pdf
- [6] Todorova, A., Syrakov, D., Gadjhev, G., Georgiev, G., Ganev, K.G., Prodanova, M., Miloshev, N., Spiridonov, V., Bogatchev, A., Slavov, K. (2010) Grid computing for atmospheric composition studies in Bulgaria. *Earth Sc. Inf.*, 3 (4), 259-282
- [7] Visschedijk, A., Zandveld P., van der Gon, H. (2007) A high resolution gridded European emission database for the EU integrated project GEMS, TNO report 2007-A-R0233/B, The Netherlands Brunekreef B, Holgate S: Air pollution and health., *Lancet* 2002, 360:1233-1242.

Performance Analysis of Real-time Applications for Debugging Parametrization

I. Georgiev, I. Georgiev

Development of safety-critical real-time applications requires from the programmers to explicitly calculating the execution time scenarios. In this paper, we perform some experimental analysis and formalization of the parameters that have significant influence to the performance of safety-critical real-time applications. Parameter's range is implemented in both hardware and software debugger that can lead the programmer to estimate the worst-case execution time (WCET).

We extract *the performance parameters* by starting from a simple computer model and adding one after another different hardware and software features.

For *simple pipelined CPU with registers only* the clocks per instruction CPI is influenced by several parameters. Most embedded systems fetch and execute an instruction in 5 stages, every stage is synchronized by a clock, the pipeline executes 5 instructions at different stages. Then the CPI is 1 (5 instructions in the pipeline, 5 clocks per instruction).

For *simple pipelined CPU considering hazards* the ideal pipelining is not possible. There are data hazards and control hazards that influence the execution time. Data hazards depend in the data dependencies in the program and are 5-8%. Every data hazard stall the pipeline for 2-4 clocks. The killers of performance are the control hazards that are caused by the branch instructions. Without prediction every branch instruction can stall the pipeline for 3 clocks. In average, every program has 25-35% branch instructions. Then the worst-case CPI both for data and control hazards will be 1.95. The CPU is 2 times slower because of the hazards.

For *simple pipelined CPU considering hazards and loop programming* the data dependencies, especially loop-carried, can strongly reduce the performance of real-time program. The execution time can be improved by loop unrolling. In a simple example of parallel loop without loop-carried dependences we show that three times loop unrolling can increase the execution time of the loop 2.7 times.

If the loop is not parallel because of the loop-carried dependences it cannot be proceed in multi-core CPU and there is no performance improvement.

Memory hierarchy utilization is an important requirement to achieve high performance. The primary memory consists of processor registers, cache level one (separate instruction and data cache), cache level two, cache level three and main memory. The exchange of the information between those levels is controlled by a memory management unit and is hidden from the OS and run-time environment. The processor issues main memory addresses, the memory management unit checks whether the data is uploaded in cache one and down through the hierarchy. If the data is in the upper level, it is a hit. If it is not, it is called miss and the information has to be uploaded from the lower level. The time to upload the missing data block from the lower level is miss penalty, when the current thread waits several hundred clock cycles.

Further, the performance analysis becomes more complicated considering *high per-*

formance array programming. Array streaming or large array files put additional requirement to the performance estimation. During the processing, the arrays are uploaded blocks by blocks from the virtual memory (if it exists) to the cache memory by the memory management hardware-software tools (memory management unit and operating system). For our further consideration let us consider only two levels in the memory hierarchy (upper level and lower memory level). If the information is not in the upper level memory, there is a cache miss or virtual memory fault (let us use a common name fault). The memory management blocks the program under execution and uploads the page or segment from the lower to the upper level memory. After that, the program is unblocked and put in the ready queue of the OS scheduler. The fault could take thousands of clock cycles. In high performance programming, it is recommended general estimation of the code performance especially of cache misses and virtual memory faults in those parts of the programs that manipulate huge information.

Let us consider nested loops that access array data stored in memory in row or column order. Programming the right nesting of the loops can make the loop expressions retrieve the data in the order in which they are stored. In case the array cannot fit in the cache or main memory pages assigned for the program, precise nesting reduces cache misses and/or virtual memory faults using all data in a cache or page before they are replaced.

For *preemptive multi-threading mode* of real time embedded system the performance depends strongly on the scheduling of processes and threads. The scheduling is performed usually by the Real Time Operating system (RTOS) and the scheduling time consists of several components:

- a. Response time, which is the time from interrupt signal to the beginning of the interrupt service by the interrupt handler;
- b. Interrupt handler time, which is the time for interrupt processing;
- c. Dispatching time of RTOS to select the next thread to run in the CPU.

Most of the RTOS declare that the time of all those three component is a predictable constant, i.e. the big-O complexity is $O(1)$. The dispatching algorithms are priority based. Priorities can be:

- a. calculated in advance;
- b. shorter threads have higher priority (Rate Monotonic Scheduling);
- c. thread with earliest (absolute) deadline has highest priority.

The execution of every thread can be preempted by a thread of higher priority in most of the embedded systems (enable and disable interrupt directives are not allowed). Performance is difficult to calculate. We suggest the influence of the unpredictable preemption to be expressed increasing the value of instruction count by some percentage, which can be experimentally calculated for the concrete configuration.

For *shared memory multiprocessing (SMP)* performance evaluation depends on several factors, most of them depend both on the environment, but also on running application and the number of the sensors and actuators that can cause additional delays.

In Single Program Multiple Data (SPMD) the same threads run in concurrent (if

the threads number is greater than the number of cores) and partly parallel mode on different cores (if the number of cores is enough). In concurrent mode the worst case is when all threads are executed in one core sequentially. In fully parallel mode, some additional delay is generated by the barrier synchronization between the threads (every thread waits for others to reach the barrier instruction).

In the *implemented experimental debugger parametrization* programming and debugging of the real time applications is supported by powerful hardware-software integrated development environment (IDE). The IDE configuration consists of:

- a. Host computer that runs programming automation package (compiler, editor, linker) and a host-debugger that runs on a simulator of the selected micro-controller;
- b. Evaluation board that has usually two micro-controllers: the target micro-controller and a microcontroller-debugger.

Both host-and-microcontroller debuggers provide a suite of functions that support the programming cycle but not so many for performance estimation. The host-debugger provides the number of the executed instructions and the simulated execution time. The microcontroller-debugger uses the incorporated timers and provide the absolute time of execution. Both simulated and absolute execution times are only snapshots of the current run of the thread. Obviously such estimation has to be modify with parameters according to the specifics of the computing configuration and possible maximum number of the connected sensors and actuators.

The proposed parametrization of the debugger has two goals. The first one is to offer the designer some percentage over the measured time. For example, statistics shows that for single ARM-based embedded systems the cache misses are 5% from all memory accesses and the miss penalty is about 5000 CPU cycles (one cycle is 5 clocks). The second goal is to incorporate entry points for some implemented algorithms for formal estimation of the worst case (there are several available open source methods). Our performance analysis has resulted in some parametrization of the microcontroller-debugger. We added some preprocessor to the debugger that performs a dialog with the designer to fill in some values for every parameter. The values can be:

- a. just Yes/No for some architectural feature of the embedded system or RTOS;
- b. some hard-coded or given by some expertize percentage to increase the WCET ;
- c. call entry point to some additional software to calculate performance. We grouped the parameters into two groups: parameters that increase the clock per instructions CPI and parameters that increase the Instruction count IC.

Below we enumerate only main groups of parameters - every group can include a lot of parameters.

The group of parameters that increase the CPI:

a. *Data and control hazards*. Here we have several entry point for subroutines that can statically or dynamically estimate the number of the data hazards and give some percentage to increase the execution time. Some example percentage is 5-8%.

b. *Branch predictions*. Several parameter are to be declared: forward branches are predicted taken, backward branches are predicted taken, dynamical branch tables organization, etc.

c. *Data dependencies especially, loop-carried dependencies*. The most important here

is the static estimation of the data dependencies - here there are entry points to different subroutine for loop-carried dependencies.

d. Single core or SMP organization. If CPU is a single core the estimation of execution time can be measurement based. In case of SMP several parameters have to be given that make the prediction more appropriate. The literature mostly considers methods for single core processors. But fully preemptive scheduling, the possible line of waiting interrupts, conflicting access to main memory makes the timing analysis dependable by the interference delays parametrization.

The parameters that increase the IC are related to:

a. Cache misses. The influence of the cache misses depends on the CPU organization, the levels of caches, the size of the cache blocks, etc. The performance prediction is based on the measurement-based timing estimation and on the experience of the designers.

b. Write policy (write through, write back). The threads that belong to the same process share the same address space and every write can invalidate some blocks. For that situation we incorporate several parameters that can make the performance prediction more correct.

c. Block replacement. The selected policy parameters (FIFO, LRU or randomly) has to be object of different estimation.

d. Memory hierarchy faults. The parameters that the designer has to present to debugger particularization together with measurement has to motivate possible estimation of the execution time, which could be more than ten times greater then the measured.

e. The parallel execution in different cores with multi-threaded synchronization by barrier instructions. Barrier synchronization in SPMD (Single Program Multiple Data) between the parallel threads (the threads are parallel but they are executed in concurrent mode) can increase the WCET with 1-3%.

f. The interrupt sequence and especially the scheduling procedure has the most significant influence on WCET. The designer present some parameters that describe more precisely the scheduling and the response time after interrupts.

g. Reset-restart for fault recovery is a popular techniques especially in periodical real-time applications. Every critical fault triggers restart of all safety-critical application that could be corrupted by the fault are executed again. We consider such restart a function that can increase the IC by 0.23% for periodical applications.

With parametrization the debugger can give more realistic picture during estimation of the execution time. Some parameters are also entry point to some additional sub-routines for performance prediction.

Computer Simulations of PM Concentrations Climate for Bulgaria

I. Georgieva, N. Miloshev

Abstract

The numerical simulations of the Particulate matters (PM) fields in Bulgaria have been performed using the US EPA Model-3 system as a modelling tool for 3D simulations and the system nesting capabilities were applied for downscaling the simulations to a 9 km resolution over Bulgaria. The national emission inventory was used as an emission input for Bulgaria, while outside the country the emissions are from the TNO high resolution inventory. The air pollution pattern is formed as a result of interaction of different processes, so if know the contribution of each, for different meteorological conditions and given emission spatial configuration and temporal behavior is very helpful for understanding the atmospheric composition and pollutants behavior. The “Integrated Process Rate Analysis” model option was applied to discriminate the role of different dynamic/chemical processes for the air pollution formation. Numerically obtained PM concentration fields of as well as of determining the contribution of different processes to the formation of surface PM concentrations will be demonstrated in the paper.

Introduction

The air is the living environment of human beings and atmospheric parameters have a great importance for the quality of life. According to the World Health Organization (WHO), between 2.5 and 11% of the total number of annual deaths are due to air pollution [10, 11]. Special attention is paid to primary emitted or secondary formed Particulate Matter (PM), which size varies from 0.01 μm to 50 μm . The particulates are separate in several fractions: PM10 (diameter <10 μm), PM2.5 (diameter <2.5 μm) and ultra-fine PM with diameter <0.1 μm (PM01). The topic is especially relevant for Bulgaria, where the situation is especially severe regarding of PM concentrations, and several times exceeded the limit values. The objective of the work is to demonstrate the numerically obtained PM concentration fields of as well as of determining the contribution of different processes to the formation of surface PM concentrations using modeling tools.

1 Methodology

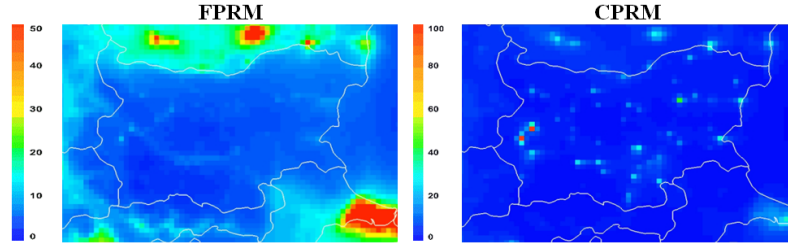
Extensive numerical simulations of the atmospheric composition fields in Bulgaria have been performed using up to date modelling tools and detailed and reliable input data [3, 4, 5]. An ensemble, comprehensive enough as to provide statistically reliable assessment of the atmospheric composition climate, has been constructed. The

used modeling tools is US EPA (Environmental Protection Agency) Model-3 system consists of: Meteorological model WRF (Weather Research and Forecasting) [7], Atmosphere Composition Model CMAQ (Community Multiscale Air Quality) [1] and Emission model SMOKE [2] (Sparse Matrix Operator Kernel Emissions). The simulations were performed day by day for 7 years (2008—2014). The NCEP Global Analysis Data meteorological background with $1^\circ \times 1^\circ$ resolution was used. The models nesting capabilities were applied to downscale the simulations to 9 km for Bulgaria. The TNO high resolution inventory [9] was exploited, and National emission inventory as input for Bulgaria. By averaging the surface concentrations over the whole simulated fields of ensemble were obtained the mean annual and seasonal surface concentrations and used as “typical” daily concentration patterns. In this work the PM are separated in 2 fractions: Fine PM (FPRM) with diameter $< 2.5\mu\text{m}$ and Coarse PM (CPRM) with diameter from $2.5\mu\text{m}$ to $10\mu\text{m}$. The Integrated Process Rate Analysis option was applied to evaluate the concentration change (ΔC) for each compound for an hour, so is presented as a sum of the contribution of the processes. The processes are advection (horizontal HADV and vertical VADV), diffusion (horizontal HDIF and vertical VDIF), emissions (EMIS), dry deposition (DDEP), chemistry (CHEM), aerosol processes (AERO) and cloud processes/aqueous chemistry (CLDS). The models computer resource requirements are rather big [8] and the numerical experiments were organized in effective HPC environment. The calculations were implemented on the Supercomputer System Avitohol at Institute of Information and Communication Technologies Bulgarian Academy of Sciences.

2 Results

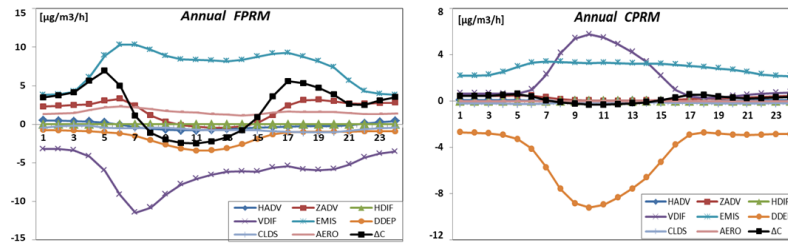
The PM climate and behavior over Bulgaria is evaluated by averaging the surface concentrations over the whole ensemble and the mean annual and seasonal surface concentrations were obtained. Due to volume limits here are present only the annually average surface concentrations for both fractions FPRM and CPRM. According to current Regulation [6] the defined limit values for PM concentrations are: FPRM – $40\mu\text{g}/\text{m}^3$ 24 hour average and $20\mu\text{g}/\text{m}^3$ annual average; CPRM – $50\mu\text{g}/\text{m}^3$ 24 hour average and $40\mu\text{g}/\text{m}^3$ annual average. The results show that there is exceedance the limit values for both PM fractions. For CPRM plots the exceedance is several times mostly at biggest cities in the country Figure 1. The outputs from the Integrated Process Rate Analysis were averaged over the 7 year ensemble and so the “typical” seasonal and annual evaluations were obtained. An example of the diurnal annual behavior of the contribution of different processes to the surface concentration of FPRM and CPRM, averaged for Bulgaria, is given in Figure 2. The processes that were considered are HADV and VADV, HDIF and VDIF, EMIS, DDEP, AERO and CLDS. The graphics show dominant contributions with their sign and phases of each process that leading to concentration change. For the FPRM can be seen that the leading processes are EMIS with positive contribution and VDIF with negative contribution. The EMIS is dominant process for CPRM too, but VDIF has highest positive contribution, and also can be see that DDEP has maximal negative contribution.

Figure 1: Annual surface concentrations for FPRM and CPRM [$\mu\text{g}/\text{m}^3$], averaged for the territory of Bulgaria at 07:00 GMT.



The ΔC has different sign during the day and depending on weather conditions and topography.

Figure 2: Annually averaged contribution of the different processes to the formation of FPRM and CPRM [$\mu\text{g}/\text{m}^3/\text{h}$] for Bulgaria.



3 Conclusion

Due to volume limitations the spatial and seasonal variability of the PM characteristics is not demonstrated at all, only the annual. The numerically obtained PM concentration fields show the several times exceeded the limit values for concentrations mainly at biggest city in the country. For the whole domain, the ΔC , leading to a change in a concentration is determined by a small number of dominating processes with big values, and the sign and phases of these processes could be opposite. The contributions sign of some processes is obvious, but for some the sign may be changing and can be different, depending on weather conditions, topography and etc.

Acknowledgment

The present work is supported by: the Bulgarian National Science Fund (grant DCVP-04/2/13.12.2016), EU -H2020 project 675121(project VI-SEEM), Program for career development of young scientists, BAS. Deep gratitude to US EPA and US NCEP, EMEP and TNO for providing free-of-charge data and software.

References

- [1] CMAQ user guide, (2006). [Online] Available from: <https://www.cmascenter.org/help/documentation.cfm?model=cmaq&version=4.6>
- [2] CEP (2003) Sparse Matrix Operator Kernel Emission (SMOKE) Modeling System, University of Carolina, Research Triangle Park, North Carolina.
- [3] Gadzhhev, G., Ganev, K., Prodanova, M., Syrakov, D., Atanasov, E., Miloshev, N. (2013) Multi-scale Atmospheric Composition Modelling for Bulgaria. NATO Science for Peace and Security Series C: Envir. Security, 137, pp. 381-385.
- [4] Gadzhhev G., K. Ganev, N. Miloshev, D. Syrakov, and M. Prodanova (2014): Analysis of the Processes Which Form the Air Pollution Pattern over Bulgaria, LNCS 8353, Springer-Verlag Berlin Heidelberg 2014, pp. 390-396
- [5] Gadzhhev G., K. Ganev, N. Miloshev, D. Syrakov, and M. Prodanova (2014): Some Basic Facts About the Atmospheric Composition in Bulgaria — Grid Computing Simulations, LNCS 8353, Springer-Verlag Berlin Heidelberg 2014, pp. 484-490.
- [6] Regulation No. 9 (State Gazette No. 46/1999, amended and supplemented, SG No. 86/2005) <http://eea.government.bg/en/output/daily/pollutants/pm.html>
- [7] Shamarock W., Klemp J., Dudhia J., Gill D., Barker D., Wang W., Powers J., (2007) A description of the Advanced Research WRF Version 2, http://www.mmm.ucar.edu/wrf/users/docs/arw_v2.pdf
- [8] Todorova, A., Syrakov, D., Gadzhhev, G., Georgiev, G., Ganev, K.G., Prodanova, M., Miloshev, N., Spiridonov, V., Bogatchev, A., Slavov, K. (2010) Grid computing for atmospheric composition studies in Bulgaria. *Earth Sc. Inf.*, 3 (4), 259-282
- [9] Visschedijk, A., Zandveld P., van der Gon, H. (2007) A high resolution gridded European emission database for the EU integrated project GEMS, TNO report 2007-A-R0233/B, The Netherlands Brunekreef B, Holgate S: Air pollution and health., *Lancet* 2002, 360:1233-1242.
- [10] World Health Organization (WHO), 2000, Fact Sheet Number 187
- [11] World Health Organization (WHO), 2004, Health Aspects of Air Pollution. Results from the WHO Project Systematic Review of Health Aspects of Air Pollution in Europe.

Numerical Methods for Fractional-in-Space Diffusion Problems

S. Harizanov, N. Kosturski, R. Lazarov, S. Margenov, P.
Marinov, Y. Vutov

1 Introduction

The interest in fractional diffusion models is motivated by the active on-going development in fractional calculus and its numerous applications related to anomalous diffusion, e.g., underground flow, diffusion in fractal domains, dynamics of protein molecules, and heat conduction with memory, to name just a few (see, e.g. [4] and the references there in). Another field of application of fractional diffusion operators is the image processing. For instance, problems with fractional power of graph Laplacians appear in image segmentation [6].

In the stationary elliptic case, such kind of problems lead to fractional order partial differential equations that involve in general non-symmetric operators. An important subclass of this topic is defined by fractional powers of self-adjoint elliptic operators, which are nonlocal but self-adjoint. In particular, the fractional Laplacian describes an unusual diffusion process associated with random excursions. In general, the fractional elliptic operators of power $\alpha \in (0, 1)$ are related to super-diffusion.

2 The problem

In what follows we consider the definition of fractional diffusion problem based on spectral decomposition of the elliptic operator. Let us consider the weak formulation of the boundary value problem: find $u \in V$ such that

$$a(u, v) := \int_{\Omega} (\mathbf{a}(x) \nabla u(x) \cdot \nabla v(x) + q(x)) dx = \int_{\Omega} f(x) v(x) dx, \quad \forall v \in V,$$

where $V := \{v \in H^1(\Omega) : v(x) = 0 \text{ on } \Gamma_D\}$, $\Gamma = \partial\Omega$, and $\Gamma = \bar{\Gamma}_D \cup \bar{\Gamma}_N$. We assume that Γ_D has positive measure, $q(x) \geq 0$ in Ω , and $\mathbf{a}(x)$ is an SPD $d \times d$ matrix, uniformly bounded in Ω . Then, the nonlocal operator \mathcal{L}^α , $0 < \alpha < 1$ is introduced through its spectral decomposition, i.e.

$$\mathcal{L}^\alpha u(x) = \sum_{i=1}^{\infty} \lambda_i^\alpha c_i \psi_i(x), \quad u(x) = \sum_{i=1}^{\infty} c_i \psi_i(x),$$

where $\{\psi_i(x)\}_{i=1}^{\infty}$ are the eigenfunctions of \mathcal{L} , orthonormal in L_2 -inner product and $\{\lambda_i\}_{i=1}^{\infty}$ are the corresponding positive real eigenvalues. Similar definition of the fractional power of a given SPD matrix is assumed.

In a very general setting, the numerical solution of nonlocal problems is rather expensive. The following four approaches A1-A4 lead to transformation of the original problem $\mathcal{L}^\alpha u = f$ to some auxiliary local problems in a computational domain of higher $(d + 1)$ dimension:

- A1** Extension to a mixed boundary value problem in the semi-infinite cylinder [2]. Truncation analysis is provided to allow numerical solution in a bounded cylinder.
- A2** Transformation to a pseudo-parabolic problem [8]. Stability conditions are obtained for the fully discrete schemes under consideration.
- A3** Integral representation of the solution in $(0, \infty)$ is used in [1]. Then exponentially convergent quadrature formulae are applied to evaluate numerically the related integrals.
- A4** Best uniform rational approximation (BURA) is introduced in [4]. See for some more details the next section.

All A1-A4 are applicable to fractional diffusion problems in computational domains with general geometry.

3 BURA methods

A class of efficient solvers for the linear system $\mathcal{A}^\alpha \mathbf{u} = \mathbf{f}$, $0 < \alpha < 1$, is proposed in [4], where \mathcal{A} is a normalized symmetric and positive definite (SPD) matrix generated by finite element or finite difference approximation of some self-adjoint elliptic problem. Instead of the original problem, the system $\mathcal{A}^{\alpha-\beta} \mathbf{u} = \mathcal{A}^{-\beta} \mathbf{f} := \mathbf{F}$, $\beta \geq 1$ an integer, is considered. Then $\mathcal{A}^{\beta-\alpha} \mathbf{F}$ is approximated by a set of solutions of systems with $\mathcal{A} + d_j \mathcal{I}$, $d_j \geq 0$, for $j = 1, \dots, k$, where $k \geq 1$ is the number of partial fractions of BURA $r_\alpha^\beta(t)$ of $t^{\beta-\alpha}$, $t \in (0, 1]$.

From algorithmic point of view, the methods A3-A4 are very similar. In both cases the approximate solution is obtained by solving a number of local problems where the sparse SPD matrix \mathcal{A} is diagonally perturbed with a positive constant. The computational efficiency follows from the assumption that some optimal (say multigrid or multilevel) solver is applied for the related auxiliary sparse SPD problems.

Both A3-A4 lead to positive approximation of the fractional diffusion operator [5]. This means that the numerical solution has monotone behaviour. This is an important advantage allowing to avoid possible numerical oscillations in the boundary layers.

When compare A3 and A4, some advantages of the BURA methods can be observed. They are stronger expressed for stronger super-diffusion, means for smaller α . A test problem with checkerboard right hand side is introduced in [1] and is used for the comparative analysis in [4]. Numerical results for the relative ℓ_2 accuracy of 1-BURA ($\beta = 1$) are presented, where the mesh parameter $h = 2^{10} \approx 10^3$ and $k = \{9, 8, 7\}$ for $\alpha = \{0.25, 0.5, 0.75\}$, respectively. For instance, for $\alpha = 0.25$, 40 is the smallest number of linear systems to be solved according to A3 to beat the numerical accuracy of 1-BURA with 10 similar systems.

4 Parallel algorithms

The supercomputing simulations are de facto a standard in the development of many new technologies, advanced engineering solutions and research projects. In particular, the use of parallel algorithms makes the fractional diffusion modeling more feasible and attractive. However, the efficient parallel computations require development of appropriate parallel algorithms. The selection of best fitted algorithms needs extended scalability analysis varying the most advanced parallel architectures.

First scalability study of parallel algorithms for methods A2-A3 are presented in [3]. Not surprisingly, more promising strong scalability results are reported for the second method. The related parallel algorithm is based on a two-level parallelization template. At the first level, a number of independent local elliptic subproblems are solved, while at the second one, parallel multigrid solvers are employed.

First robustness and parallel scalability results for BURA are published in [4]. A more involved performance analysis of multigrid preconditioning (utilized in the related BURA implementation) on Intel Xeon Phi towards scalability for extreme scale problems is available in [7].

The parallel numerical tests included in the above mentioned papers are run on HPC cluster Avitohol at ICT-BAS.

5 Concluding remarks

Some recently introduced numerical methods for fractional-in-space diffusion problems are discussed with a particular focus on the method based on on best uniform rational approximation (BURA) of $t^{\beta-\alpha}$ for $0 \leq t \leq 1$ and natural β . Bigger β means stronger regularity assumptions. This is why, from practical point of view, $\beta = 1$ is the most important case. Unlike the integral quadrature formula method (A3) the approximation properties of the BURA algorithm are not symmetric with respect to $\alpha = 0.5$, $\alpha \in (0, 1)$. Some well expressed advantages for smaller α , i.e., in the case of stronger super-diffusion, are observed.

The plans for future research include: i) developing BURA of higher order which requires also improvements of the involved Remez algorithm; ii) developing new approaches based on BURA including novel multi-step methods; iii) developing BURA based methods for more general or/and new nonlocal problems.

Acknowledgement

This research has been partially supported by the Bulgarian National Science Fund under grant No. BNSF-DN12/1. The work of S. Harizanov and Y. Vutov has been also partially under grant No. BNSF-DM02/2. The work of R. Lazarov is partially supported by Grant NSF-DMS 1620318.

For running the parallel numerical tests, we acknowledge the access to the HPC cluster Avitohol at ICT-BAS.

References

- [1] A. Bonito, J. Pasciak, *Numerical Approximation of Fractional Powers of Elliptic Operators*, Mathematics of Computation, 84 (2015), 2083-2110
- [2] L. Chen, R. Nocketto, O. Enrique, A.J. Salgado, *Multilevel Methods for Nonuniformly Elliptic Operators and Fractional Diffusion*, Mathematics of Computation, 85 (2016), 2583-2607
- [3] R. Ciegis, V. Starikovicius, S. Margenov, R. Kriauzien, *Parallel Solvers for Fractional Power Diffusion Problems*, Concurrency Computat: Pract Exper (2017), e4216, <https://doi.org/10.1002/cpe.4216>
- [4] S. Harizanov, R. Lazarov, S. Margenov, P. Marinov, Y. Vutov, *Optimal Solvers for Linear Systems with Fractional Powers of Sparse SPD Matrices*, Numerical Linear Algebra with Applications (2018) DOI: 10.1002/nla.2167
- [5] S. Harizanov, S. Margenov, *Positive Approximations of the Inverse of Fractional Powers of SPD M-Matrices*, to appear in Springer Lecture Notes in Economics and Mathematical Systems, available as arXiv:1706.07620, 2018
- [6] S. Harizanov, S. Margenov, P. Marinov, Y. Vutov, *Volume constrained 2-phase segmentation method for utilizing a linear system solver based on the best uniform polynomial approximation of $x^{-1/2}$* , Journal of Computational and Applied Mathematics, 310 (2017), 115-128
- [7] N. Kosturski, S. Margenov, Y. Vutov, *Performance Analysis of MG Preconditioning on Intel Xeon Phi: Towards Scalability for Extreme Scale Problems with Fractional Laplacians*, Large-Scale Scientific Computing, Springer LNCS, Vol. 10665 (2018), 304-312
- [8] P.N. Vabishchevich, *Numerically Solving an Equation for Fractional Powers of Elliptic Operators*, Journal of Computational Physics, 282 (2015), 289-302

Edge Detection of Radiographic Images through Phantom Blur Denoising

S. Harizanov, I. Lirkov, I. Georgiev, J. Stary, S. Zolotarev

1 Theoretical setup

Digital image processing is a vastly emerging and extremely active research field, that combines mathematical and computer science knowledge, and has applications in practically all our daily activities. This paper is devoted to a novel approach for extracting important structural information for the scanned object, based on the segmented difference image between two denoised outputs. More precisely, denoting by $\bar{\mathbf{u}}$ the noise-free “perfect” gray-scale image we want to reconstruct, and assuming that we have two independent noisy observations \mathbf{f}^1 and \mathbf{f}^2 of it with identical noise characteristics, we can study the following two model problems

$$\mathbf{u}_{UB} = \operatorname{argmin}_{u \in [0, \nu]^n} \|\nabla u\|_{2,1} \text{ s.t. } \begin{cases} \|T(u) - T(\mathbf{f}^1)\|_2^2 \leq \frac{1}{2} \|T(\mathbf{f}^2) - T(\mathbf{f}^1)\|_2^2 \\ \|T(u) - T(\mathbf{f}^2)\|_2^2 \leq \frac{1}{2} \|T(\mathbf{f}^2) - T(\mathbf{f}^1)\|_2^2 \end{cases}, \quad (1)$$

and

$$\mathbf{u}_B = \operatorname{argmin}_{u \in [0, \nu]^n} \|\nabla u\|_{2,1} \text{ s.t. } \begin{cases} \|T(Hu) - T(\mathbf{f}^1)\|_2^2 \leq \frac{1}{2} \|T(\mathbf{f}^2) - T(\mathbf{f}^1)\|_2^2 \\ \|T(Hu) - T(\mathbf{f}^2)\|_2^2 \leq \frac{1}{2} \|T(\mathbf{f}^2) - T(\mathbf{f}^1)\|_2^2 \end{cases}. \quad (2)$$

Here, n is the image size, ν is its maximal intensity, T is the *Anscombe transform*, $\nabla \in \mathbb{R}^{2n \times n}$ is the discrete gradient operator (forward differences and Neumann boundary conditions are used), the Total Variation (TV) norm $\|\cdot\|_{2,1}$ sums the lengths of the pixel gradients, and $H \in [0, \nu]^{n \times n}$ is a blur operator, corresponding to a convolution with a Gaussian kernel of standart deviation $\sigma = 0.5$.

The constrained energy minimization techniques (1) and (2) lead to practically noise-free, but oversmoothed images. The level of oversmoothing depends on the characteristics of the constrained set and the admissibility of the true image $\bar{\mathbf{u}}$. It affects the contrast of the image edges and not the brightness of regular regions. Therefore, we can try localizing the image edges via looking in the segmented difference image

$$|\mathbf{u}_{UB} - \mathbf{u}_B| > \text{threshold}.$$

2 Numerical results

The work of S. Harizanov is supported by the Bulgarian National Science Fund under grant No. BNSF-DM02/2.

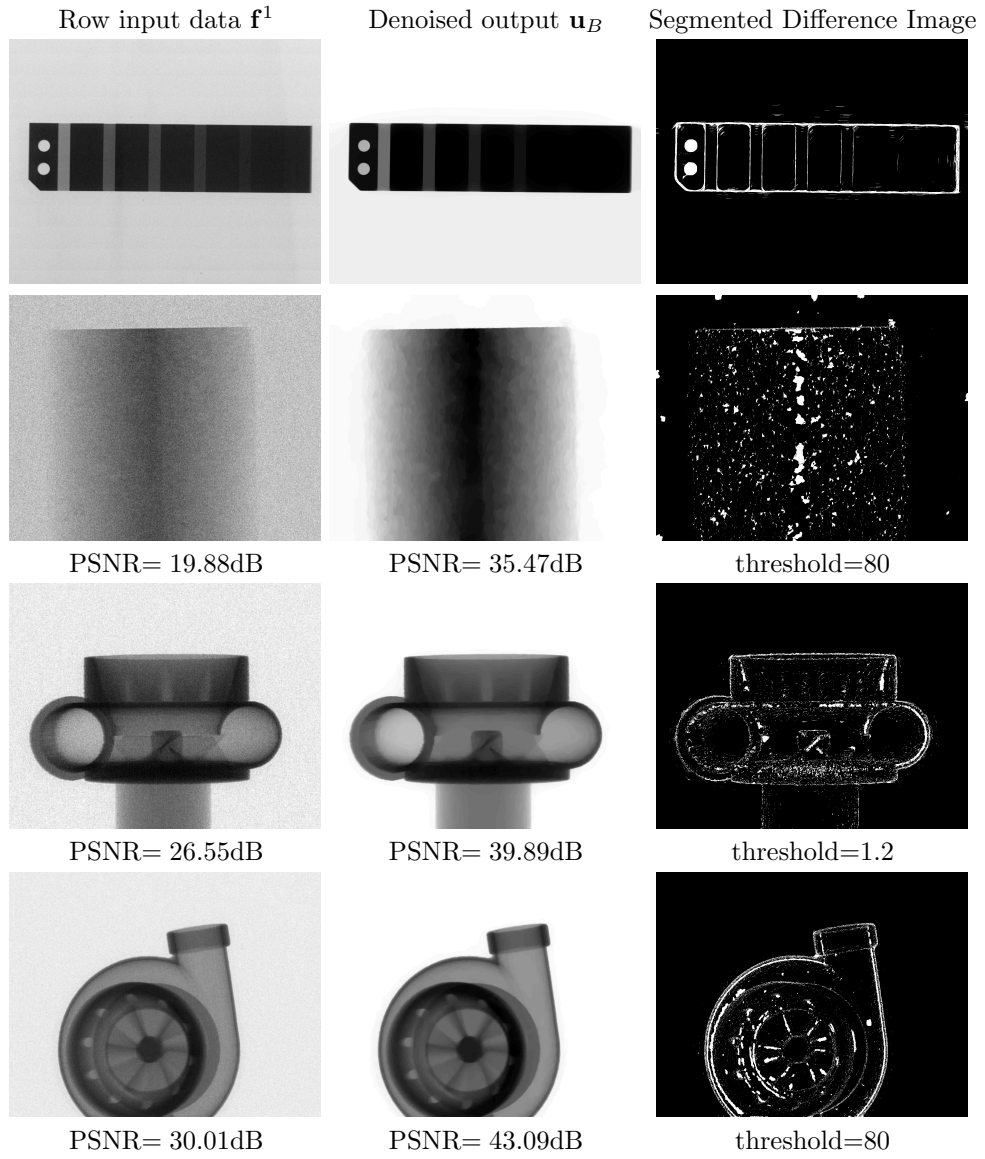


Figure 1: Experimental verification of the proposed methodology on various data sets, generated by the industrial tomograph Nikon XT H 225.

Study of Non-Proline *cis* Peptide Planes in Different Protein Framings

Y. Hou, J. Dai, A. J. Niemi, X. Peng, J. He, N. Ilieva

1 Introduction

Familiar examples where the choice of coordinates is important, range from description of Keplerian planetary motion to identification of action-angle variables in integrable models [1]. In studying and visualisation of proteins, however, this problem has been widely neglected so far. Commonly, a protein is visualised as a one-dimensional piecewise linear discrete chain with vertices that coincide with the positions of backbone C^α atoms. Geometrical description of such an object can be performed both in intrinsic and extrinsic coordinates. A classical differential-geometry example from the first category are the canonical Frenet coordinates [2]. More common, however, is the description in terms of Ramachandran angles — the dihedral angles that are adjacent to the α -carbons of the backbone [3]. However, most of the 3D visualisation programs, e.g. VMD, Jmol and PyMOL [4], employ an extrinsic, laboratory coordinate frame. The extrinsic and intrinsic geometries are the same for structureless curves, but their information content differs for natively framed objects as is the case with proteins. We step on the *intrinsic* geometric structure provided by the peptide planes to develop a methodology for analysis and visualisation of protein structure with the potential for identifying the intrinsic geometry influence on atomic positions anomalies and overcoming the apparent statistical bias [5] in the diffraction data refinement.

2 Materials and Methods

For description and analysis of protein structure and dynamics we employ, on the one side, the *extrinsic* Frenet frames (see, e.g. [6]). On the other, based on the planar character of the peptide bond, we introduce an *intrinsic* coordinate system, defined by the C_i , N_i and O_i atoms of a given peptide plane (Fig. 1) with its origin at the location of the corresponding C_i atom as follows:

$$\mathbf{x}_i = \frac{\mathbf{r}_{O_i} - \mathbf{r}_{C_i}}{|\mathbf{r}_{O_i} - \mathbf{r}_{C_i}|}, \quad \mathbf{z}_i = \frac{\mathbf{x}_i \times \mathbf{u}_i}{|\mathbf{x}_i \times \mathbf{u}_i|}, \quad \mathbf{y}_i = \frac{\mathbf{z}_i \times \mathbf{x}_i}{|\mathbf{z}_i \times \mathbf{x}_i|} \quad (1)$$

where

$$\mathbf{u}_i = \frac{\mathbf{r}_{N_i} - \mathbf{r}_{C_i}}{|\mathbf{r}_{N_i} - \mathbf{r}_{C_i}|}$$

3 Results and Discussion

Our empirical data set consists of the ultrahigh-resolution (better than 1.0 Å) structures in PDB [7]. Therein, the peptide planes are mostly found in *trans* conformation,

defined through Ramachandran dihedral ω between $\langle C_i^\alpha - C_i - N_i \rangle$ and $\langle C_i - N_i - C_{i+1}^\alpha \rangle$ planes having the value $\omega \approx \pi$. The *cis* conformation, where $\omega \approx 0$, is seldom seen among the low-resolution structures, but still present among the less refinement-prone high-resolution ones.

We focus on these relatively rare *cis* peptide plane conformation, extending the definition range for the Ramachandran dihedral ω from the ideal value $\omega \approx 0$ to the interval $\omega \in [-\pi/4, \pi/4]$. To narrow further the selection, we choose to analyse the structures that exhibit *cis* X-Xnp peptide plane, where X denotes any amino acid and Xnp stands for any amino acid but proline. About 40% of ultrahigh-resolution structures presently deposited in PDB have at least one *cis* peptide plane, where the rare *cis* X-Xnp type accounts for about 7.5% of the structures. In terms of peptide planes, the rarity of these conformations becomes even more apparent: out of 118226 peptide planes present in the ultrahigh-resolution data set, only about 0.3% appear in *cis*-conformation, and only some 0.05% are of *cis* X-Xnp type. In our numbering convention, the i^{th} peptide plane connects residues (resp., α -carbons) at position i and $i + 1$. In Fig. 2, the frequency of different amino acids at these two positions for the *cis* X-Xnp peptide planes is shown.

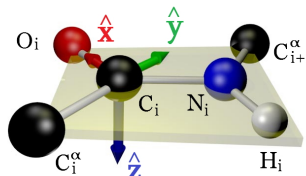


Figure 1: The right-handed orthonormal $(\mathbf{x}, \mathbf{y}, \mathbf{z})$ CNO-frame.

resolution structures presently deposited in PDB have at least one *cis* peptide plane, where the rare *cis* X-Xnp type accounts for about 7.5% of the structures. In terms of peptide planes, the rarity of these conformations becomes even more apparent: out of 118226 peptide planes present in the ultrahigh-resolution data set, only about 0.3% appear in *cis*-conformation, and only some 0.05% are of *cis* X-Xnp type. In our numbering convention, the i^{th} peptide plane connects residues (resp., α -carbons) at position i and $i + 1$. In Fig. 2, the frequency of different amino acids at these two positions for the *cis* X-Xnp peptide planes is shown.

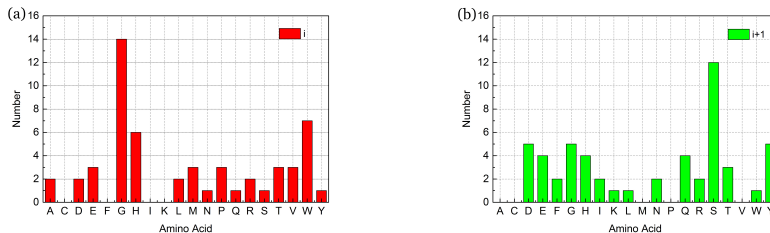


Figure 2: Distribution of amino acids at the i^{th} (left panel) and $(i + 1)^{\text{th}}$ (right panel) vertex of the *cis* X-Xnp peptide planes within ultrahigh-resolution structures in PDB.

As a characteristic example we consider the side-chain orientation in the *cis* peptide planes, substantiated by the angular distribution of the β -carbons [8] at the respective nodes of the protein backbone. In Fig. 3, Fig. 4, these distributions on a C^α -centered Frenet sphere for the i^{th} , resp. $(i + 1)^{\text{th}}$, C^β are given, in parallel for the *cis* X-Xnp and the *cis* X-Pro peptide bonds. In all plots we observe a clear deviation from the grey background of all PDB structures with resolution below 1.0 Å.

Thus, in the Frenet-representation, i^{th} *cis* peptide plane affects both the preceding and subsequent structures, as reflected in the C_i^β and C_{i+1}^β distributions, the latter though stronger. In the CNO frame (figures not shown), no influence on the preceding structures can be detected — distributions of all involved atoms are well in line with

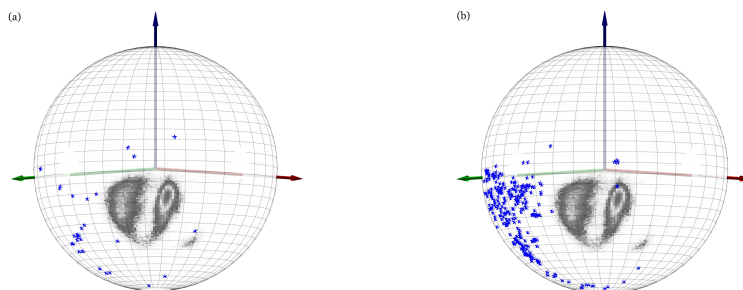


Figure 3: Distribution of the β -carbons of the *cis* X-Xnp residues (left panel) and of the *cis* X-Pro ones (right panel), on a C^α -centered Frenet sphere for the i^{th} C^β .

the corresponding *trans*-conformation background.

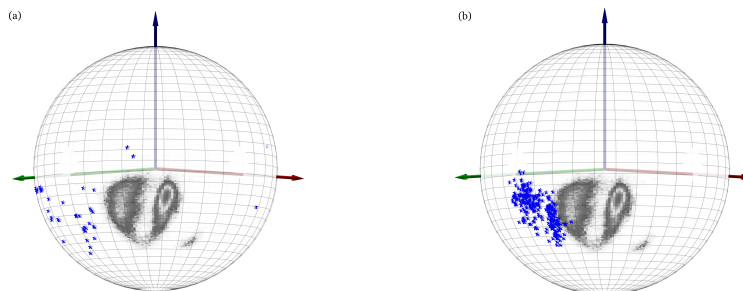


Figure 4: Same as Fig. 3, but for the $(i + 1)^{\text{th}}$ C^β .

The CNO-frame information content is complementary to that of the Ramachandran representation. However, one would expect that the CNO angular variables, being spherical coordinates, are more convenient for protein structure studies than the toroidal Ramachandran angles [9].

4 Conclusions and outlook

We compare the local geometry around *trans* and *cis* peptide planes, as well as the differences between their conformations in different coordinate systems. In particular, we show how to employ the intrinsic geometry to visually analyze the atomic level neighborhood around a peptide plane and systematically classify how the *cis* conformation affects the protein structure. Further, we combine the traditional Ramachandran angles with our modern visualisation methods. We reveal systematics in the way how such a *cis* peptide plane deforms the atomic level geometry in its neighborhood, and show how our 3D methodology easily detects the presence of a *cis* peptide planes from the arrangement of atoms near the latter. Our approach identifies

efficiently exceptionally positioned atoms in crystallographic PDB structures. It can also help overcome the apparent refinement bias towards statistically more significant *trans* structures. Thus it can be extended to a visual analysis and refinement tool, applicable even when resolution is limited or data are incomplete.

Acknowledgements

This work was supported in part by Bulgarian Science Fund (Grant DNTS-China-01/9/2014), Vetenskapsrådet (Sweden), Carl Trygger's Stiftelse and Qian Ren Grant at Beijing Institute of Technology.

References

- [1] L.D. Faddeev, L.A. Takhtajan, Hamiltonian methods in the theory of solitons (*Springer Verlag, Berlin, 1987*)
- [2] M. Spivak, *A Comprehensive Introduction to Differential Geometry Vol. 5* (Publish or Perish, Houston, 1979)
- [3] G. N. Ramachandran, C Ramakrishnan, V. Sasisekharan, *J. Mol. Biol.* **7** 95(1963)
- [4] https://en.wikipedia.org/wiki/List_of_molecular_graphics_systems
- [5] A. Jabs, M.S. Weiss, R. Hilgenfeld, *Non-proline cis peptide bonds in proteins.* *J. Mol. Biol.* **286** 291(1999)
- [6] S. Hu, M. Lundgren, A.J. Niemi, *The Discrete Frenet Frame And Curve Visualization with Applications to Folded Proteins.* *Phys. Rev.* **E83** 061908 (2011)
- [7] H.M. Berman, J. Westbrook, Z. Feng, G. Gilliland, T.N. Bhat, H. Weissig, I.N. Shindyalov, P.E. Bourne, *The protein data bank.* *Nucl. Acid Res.* **28** 235(2000)
- [8] M. Lundgren, A.J. Niemi, *Correlation between protein secondary structure, backbone bond angles, and side-chain orientations.* *Phys. Rev.* **E86** 021904 (2012)
- [9] Y. Hou, J. Dai, A.J. Niemi, X. Peng, J. He, and N. Ilieva, *Intrinsic protein geometry with application to non-proline cis peptide planes* (in preparation)

Study of Human Interferon-Gamma Glycosylation by Molecular Dynamics Simulations

E. Lilkova, N. Ilieva, P. Petkov, L. Litov

1 Introduction

Interferon- γ (IFN γ) is a signaling molecule, which is crucial in regulation of the formation and modulation of immune response. Human IFN γ (hIFN γ) consists of 143 amino acids (aa), organised in six α -helices (comprising 62% of the molecule), which are linked by short unstructured regions. Besides them, hIFN γ contains also a long positively charged unstructured C-terminal domain composed of 21 aa. The mature form of hIFN γ is organised as a non-covalent homodimer. The cytokine accomplishes its functions via high-affinity extracellular interaction with its specific receptor (IFN γ R1).

Under physiologic conditions the natural human IFN γ is a glycoprotein. The two N-glycosylation sites in each monomer chain – ASN²⁵ and ASN⁹⁷ – are independently and differentially glycosylated. Naturally, two fractions with molecular weights of 20 kDa and 25 kDa are isolated, that correspond to either monoglycosylated or diglycosylated protein. The chemical compositions of the predominantly occurring oligosaccharide sequences in the native hIFN γ have been determined experimentally [3] and are shown in Figure 1. Glycosylation does not affect hIFN γ activity, but it has been shown that N-linked oligosaccharides promote the folding and dimerization of the recombinant cytokine. In addition, N-glycosylation enhances circulatory half-life by protecting hIFN γ from proteolytic degradation [4].

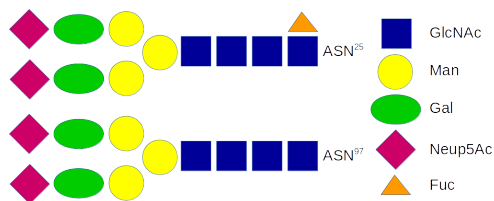


Figure 1: Chemical composition of the most common oligosaccharide chains of native hIFN γ .

Here, we report the development of model structures of monoglycosylated at either ASN²⁵ or ASN⁹⁷ or diglycosylated full-length native dimers and study the influence of added glycans on the conformation and dynamics of hIFN γ by means of molecular dynamics (MD) simulations.

2 Materials and Methods

Input model development

We started from the conformation of previously reconstructed and folded full-length hIFN γ [5]. Glycosylation anchors were added to the starting structure of the protein using GlyProt server [1]. The glycan structures, corresponding to Figure 1, were then built and attached to the anchors with the help of Glycan Reader of the molecular dynamics simulation setup server CHARMM-GUI [2]. The structures were parameterized with the CHARMM 36 force field, minimised and equilibrated.

Production simulation protocol

The production MD simulations were performed with GROMACS 2016.3. The glycoproteins were solvated in rectangular boxes with a minimal distance to the box walls of 2 nm under periodic boundary conditions. The leapfrog integrator with constraints imposed on all bonds was used with a time-step of 2 fs. Van der Waals interactions were smoothly switched off from a distance of 10 Å and truncated at 12 Å. The direct PME cut-off was 12 Å. The simulations were performed at temperature of 310 K and pressure of 1 atm, supported through a v-rescale thermostat and Parrinello-Rahman barostat, respectively.

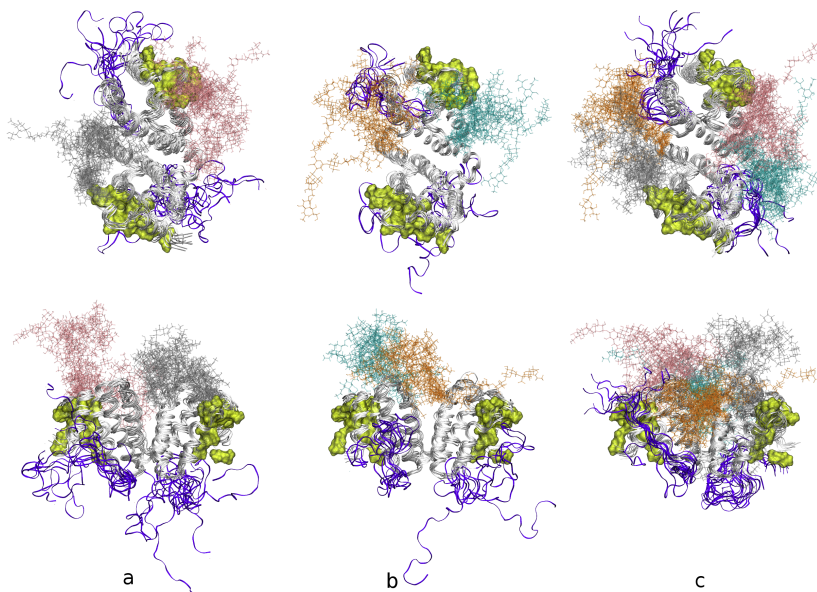


Figure 2: Top and front projections of the sampled conformations of (a) monoglycosylated at ASN²⁵, (b) monoglycosylated at ASN⁹⁷ and (c) diglycosylated full-length hIFN γ . The globular part of the proteins is depicted in grey ribbons and the flexible C-termini are in blue ribbons. The glycans are depicted in lines as follows: ASN²⁵_A – silver, ASN⁹⁷_A – cyan, ASN²⁵_B – pink, and ASN⁹⁷_B – orange.

3 Results and Discussion

The glycan chains populate conformations, which remain in the upper part of the hIFN γ molecule. A graphical summary of the three trajectories is shown in Figure 2. The carbohydrates do not alter significantly the conformation and dynamics of the globular part of hIFN γ . In fact, they even stabilise the globule and reduce its atoms RMSF (Figure 3).

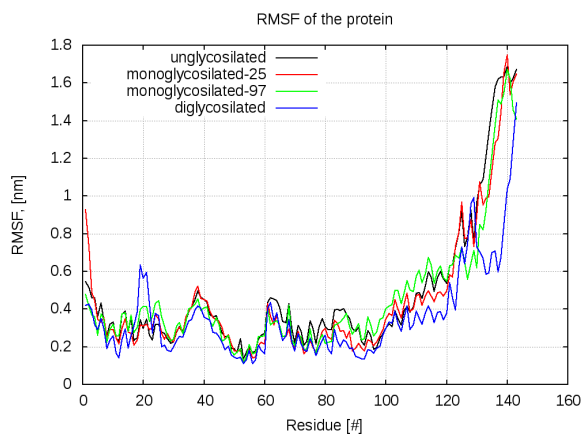


Figure 3: RMSF of the amino acid residues of hIFN γ .

The glycan chains interact both with the receptor-binding interfaces and the flexible C-termini of the cytokine. This is reflected in the contacts, formed between the carbohydrates and different parts of hIFN γ , presented in Figure 4. As evident from the two top plots, the glycans at position ASN²⁵ interact more actively with the receptor-binding sites than those at position ASN⁹⁷.

Nonetheless, contrary to expectations, the glycans do not interact very actively with the C-termini, except for the diglycosylated glycoprotein. Initially in all three simulations, the C-termini are close to the globular part. In both cases of monoglycosylated hIFN γ one or both tails drift away from the globule into the solvent, which prevents them from binding to the glycans, added to the upper part of the cytokine globule. Thus, our results do not support the hypothesis that interaction between the oligosaccharides and the C-termini protects the latter from proteolytic degradation.

Acknowledgements

The simulations were performed on the supercomputer Avitohol@BAS and on the HPC Cluster at the Faculty of Physics of Sofia University St. Kl. Ohridski. This research was supported in part under the Programme for young scientists' career development at the Bulgarian Academy of Sciences (DFNP-17-146/01.08.2017).

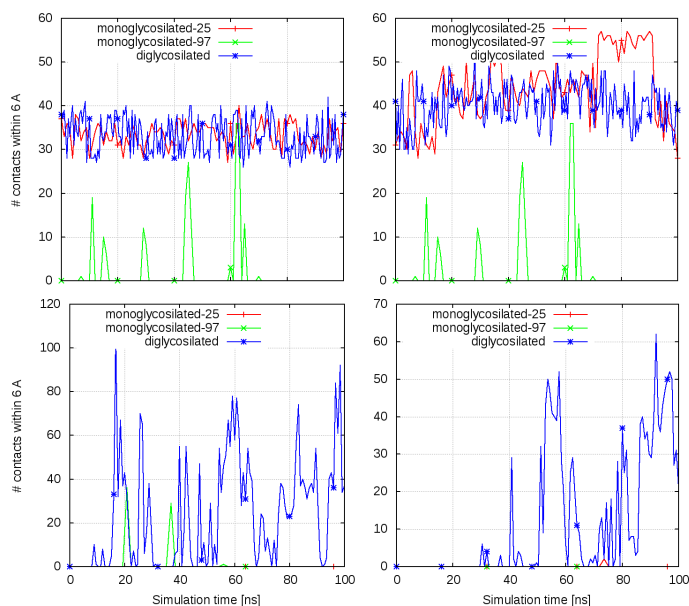


Figure 4: Number of contacts within 6 Å between the carbohydrate chains and the left and right receptor-binding interfaces (left and right top panels) and the C-termini of the two monomers (left bottom panel – chain A, right bottom panel – chain B) of hIFN γ .

References

- [1] <http://www.glycosciences.de/modeling/glyprot/php/main.php>.
- [2] <http://www.charmm-gui.org/?doc=input/glycan>.
- [3] T Sareneva, E Mortz, H Tolo, P Roepstorff, I Julkunen, Eur J Biochem (1996) 242, 191-200.
- [4] A Razaghi, L Owens, K Heimann, J Biotechnol (2016) 240, 48-60.
- [5] P Petkov, E Lilkova, N Ilieva, G Nacheva, I Ivanov, and L Litov, In: Large-Scale Scientific Computing. LSSC 2017. Lecture Notes in Computer Science, vol. 10665, (2017), 544-551.

Monte Carlo Simulation for Seismic Analysis of Egnatia Highway Bridges in Northern Greece

K. Liolios, T. Makarios, A. Liolios, K. Georgiev, I. Georgiev

1 Introduction

A probabilistic methodology concerning the construction of vulnerability curves for Civil Engineering Structures under seismic excitation, and especially for bridges, is presented. Use is made of the *Finite Element Method* (FEM), the non-linear static pushover procedure, the capacity spectrum method and Monte Carlo simulation techniques for the treatment of various uncertainties. The methodology is applied to obtain the fragility curves of the ravine *Egnatia Highway bridge* in the Kavala motorway section, East Macedonia, Greece.

The vulnerability analysis of Civil Engineering Structures, and especially for highway bridges, represents a critically important step in their seismic damage estimation and protection process [1]. The relevant fragility curves provide the probability that a specific damage level will be exceeded for a given intensity of a seismic event. In this respect, development of vulnerability relationships for both, the existing and under design Civil Engineering structures, is a key element in formulating mitigation and disaster planning strategies in Civil Earthquake Engineering for the estimation of the urban seismic risk. In combination with seismic hazard analysis at the bridge sites, they can lead to a reliable assessment of the seismic risk of highways. Furthermore, they can even be used by the authorities in charge to prioritize the on site aftershock inspections, in order to check the structural integrity of the bridges subjected to a severe seismic event.

The present paper deals with a simplified analytic methodology for the evaluation of vulnerability curves for bridges. The methodology combines the nonlinear static pushover procedure, the capacity spectrum method [2], and Monte Carlo simulation techniques for the treatment of various uncertainties [3, 4]. The methodology is applied for establishing fragility curves for an reinforced concrete bridge in the Kavala section of Egnatia Motorway, in the county of East Macedonia, Northern Greece. The Kavala bridge examined herein is a structurally representative one of many bridges in Egnatia Motorway, and in Greece more generally [5-7].

Egnatia Odos is a new motorway that crosses Northern Greece in an E-W direction. It is currently the largest and technically the most demanding highway project in Greece, and one of the biggest ones under recent (2008-2009) construction in Europe. Moreover, for the design and construction of Egnatia Motorway, a lot of Applied Science topics are involved, e.g. structural and seismic mechanics, geotechnical and transport engineering, hydraulic and environmental engineering, probabilistic methods, etc. The total length of Egnatia Motorway is about 1000 km and includes about 1900 special structures (bridges, tunnels and culverts). These structures are expected to withstand several minor or moderate earthquakes during their life, and may be

damaged if they are subjected to a major (catastrophic) earthquake. So, the construction of their fragility curves is very significant [5].

2 Problem statement and solution procedure

As was mentioned in the Introduction, the present study focuses on the simplified practical fragility analysis of bridges. Details have been presented in [7]. The vulnerability functions, required for the fragility curves, are expressed [1, 6-7] in terms of a Lognormal cumulative probability function in the form of next eq. (1):

$$P_f(DP \geq DP_i|S) = \Phi \left[\frac{1}{\beta_{tot}} \cdot \ln \left(\frac{S}{S_{mi}} \right) \right] \quad (1)$$

Here $P_f()$ is the probability of the damage parameter DP being at, or exceeding, the value DP_i for the i -th damage state for a given seismic intensity level defined by the earthquake parameter S (the Peak Ground Acceleration-PGA or Spectral Displacement - S_d), Φ is the standard cumulative probability function, S_{mi} is the median threshold value of the earthquake parameter S required to cause the i -th damage state, and β_{tot} is the total lognormal standard deviation. Thus, the description of the fragility curve involves the two parameters, S_{mi} and β_{tot} , which must be determined. The damage level depends on the input seismic excitation, i.e. the seismic ground acceleration. As well known from Structural Dynamics and Earthquake Engineering [2], because this input is not known for future earthquakes, the spectral approach is used according to various seismic building codes, e.g. the Greek Aseismic Code EAK2000 [5].

According to equation (1), the description of the fragility curve involves only two parameters, S_{mi} and β_{tot} . The first parameter S_{mi} is estimated on the basis of the capacity spectrum method [1], wherein the demand spectrum is plotted for a range of values of the earthquake parameter S (in spectral acceleration vs. spectral displacement format) and it is superimposed on the same plot with the capacity curve of the bridge. The earthquake parameter used in this study is the peak ground acceleration (PGA).

The second parameter of Eq. (1) is the total lognormal standard deviation β_{tot} , which incorporates the various uncertainties in the seismic demand, in the response and the capacity of the bridge, and also in the definition of the damage index and damage states. So, it takes into account the uncertainties in seismic input motion (demand), in the response and resistance of the bridge (capacity), and in the definition of damage states. This parameter (β_{tot}) can be estimated in the frame of Monte Carlo simulation techniques [3] by realizing a statistical combination of the individual uncertainties, assuming these are statistically independent. On the basis of empirical fragility curves obtained from actual Egnatia Highway bridges damage data in the frame of the research project ASPROGE [5], the value of β_{tot} was set equal to 0.60

3 The investigated case of the Kavala ravine bridge in Egnatia Motorway

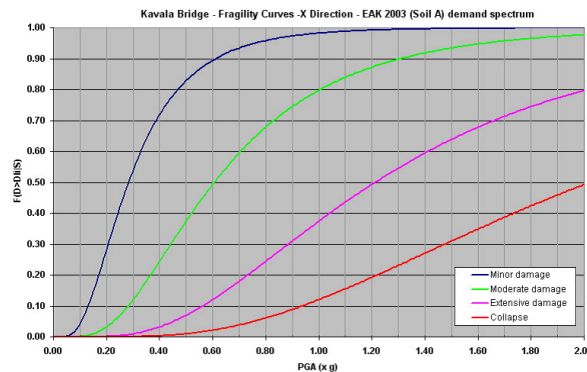
On the part of Egnatia Motorway, northern Greece, that bypasses the city of Kavala, East Macedonia, it has been constructed a ravine bridge, shown in Fig. 1. The bridge is of reinforced concrete, with seismic stoppers, and crosses a 54.00 m deep ravine. The 180 m long bridge consists of four spans, each constructed using four 45 m long prestressed beams that rest on three piers and two abutments via elastomeric bearings.

Figure 1: The Kavala ravine bridge on Egnatia Motorway.



The fragility curves were computed, assuming a lognormal cumulative probability distribution for the damage ratio as a function of peak ground acceleration PGA and using the equation (1). Representative results for the Kavala bridge obtained by the numerical implementation of the proposed methodology are presented in Fig. 2. These results concern the fragility curves in the in the x (longitudinal) direction.

Figure 2: Fragility curves for the Kavala bridge (x longitudinal - direction, EAK2003 elastic demand spectrum).



Finally, for the case of bridges with seismic stoppers, where unilateral contact effects must be taken into account, the approach presented in [7] using the hemivariational inequality concept [8] can be applied.

References

- [1] Shinozuka, M., Feng, M.Q., Lee, J., Naganuma, T., *Statistical analysis of fragility curves*, J. Eng. Mech. 126(12), 1224-1231 (2000)
- [2] Chopra, A.K., *Dynamics of Structures. Theory and Applications to Earthquake Engineering*, Pearson Prentice Hall, New Jersey (2007)
- [3] Dimov, I.T., *Monte Carlo Methods for Applied Scientists*, World Scientific Publishing Company, Singapore (2008)
- [4] Hwang, H.H.M, Jaw, J.W.: Probabilistic damage analysis of structures. J. Struct. Eng. 116(7), 1992-2007 (1990)
- [5] ASPROGE. *Research Project for the Aseismic PROtection of Bridges. Egnatia Odos S.A., Thessaloniki* (2007)
- [6] Moschonas, I.F., Kappos, A.J., Panetsos, P., Papadopoulos, V., Makarios, T., Thanopoulos, P., *Seismic fragility curves for Greek bridges: Methodology and case studies*. B. Earthq. Eng. 7(2), 439-468 (2009)
- [7] Liolios, Ast., Panetsos, P., Liolios, Ang., Hatzigeorgiou, G., Radev, S., *A numerical approach for obtaining fragility curves in seismic structural mechanics: A bridge case of Egnatia Motorway in northern Greece*, In: Dimov I., Dimova S. and Kolkovska N. (Eds.), Numerical Methods and Applications, Lecture Notes in Computer Science (LNCS) 6046, 477-485 (2010)
- [8] Panagiotopoulos, P.D., *Inequalities and Applications in Mechanics and Engineering*, Springer Verlag, Berlin (1993)

Scalability Analysis of Solvers based on Hierarchical Compression of Dense Matrices and Gaussian Elimination

D. Slavchev, S. Margenov

Abstract

We compare the performance of traditional Gaussian elimination with a solver utilizing hierarchical compression of the matrix. The test problems are obtained by Boundary Element Method (BEM) simulation of laminar flow around airfoils. The most computationally expensive part of the BEM algorithm is to solve the arising system of linear algebraic equations. The related dense matrix can be compressed using a Hierarchically Semi-Separable (HSS) representation. This significantly lowers the computational complexity of the solution method, thus allowing faster overall execution.

The performance of STRUMPACK library implementation of HSS and the MKL direct solver is compared on Intel Xeon CPU architecture. At the end, we examine the accuracy of the HSS approximation using the (exact) results of Gaussian elimination as a reference solution.

1 Background

This work is motivated by the recent development of heterogeneous high performance computing (HPC) architectures. Solving systems of linear algebraic equations with dense matrices is one of the most computationally intensive numerical linear algebra problems. This is the topic of the article. The focus is on a comparative performance analysis of a solution method based on hierarchical compression of a class of test matrices obtained by BEM simulation of laminar flows around airfoils.

The traditional Gaussian elimination has computational complexity $O(n^3)$, where n is the number of unknowns (degrees of freedom). The methods based on hierarchical compression have nearly optimal complexity, i.e., $O(r^2n)$, where r is the maximum rank of off diagonal blocks of the matrix. Typically r is much smaller than n . For some problems it is either a constant or it grows slowly with n .

The organization of this short communication is as follows. A brief overview of the HSS implementation in STRUMPACK is given in Section 2. Some numerical results are presented in Section 3 ending with a short summary in Section 4.

2 HSS method

A summary of Hierarchically Semi-Separable matrices can be seen in [1], while the parallel algorithm implemented in the STRUMPACK package is described in [2]. The HSS framework developed in STRUMPACK consists of HSS compression, ULV

factorization and solution. The HSS compression is the most important part of the HSS framework. Once the matrix is compressed specialized fast operations can be performed on it, including factorization and matrix vector product.

The HSS compression uses a *cluster tree* that defines a hierarchical partitioning of the dense matrix A . This decomposition can be performed for any matrix, but it has a practical value mostly when the off-diagonal blocks of A have a low-rank. The algorithm uses randomized sampling, which is implemented by multiplying the matrix with a set of random vectors. This method is introduced by Martinsson [1]. The main advantage is that it doesn't need to access the entire matrix A , but only parts of it. If a fast sparse matrix-vector product is utilized, the complexity of the HSS compression is $O(r^2n)$, where n is the size of the matrix and r is the maximum rank of the off diagonal blocks, found during the compression. For structured matrix, like the ones produced by the BEM method, r is much smaller than n .

In order to calculate the maximum rank r , STRUMPACK should be supplied with an user defined compression tolerance ϵ . If this tolerance is too small (e.g. approaching the machine accuracy), then r will be close to n , and the matrix will not be compressed significantly. This means to have a computational complexity close to $O(n^3)$. If the tolerance is too bigger the complexity of the HSS method will be much smaller, but at the price of accuracy of the numerical solution.

The compressed matrix is then factorized using ULV-like factorization. The special structure of the HSS matrix is taken into account. Finally the ULV factorized form is used to compute the solution.

3 Numerical results

The presented numerical results are obtained on the HPC cluster AVITOHOL of the Institute of Communication and Information Technologies, Bulgarian Academy of Sciences. We run the tests on a single node with two Intel Xeon E5-2650v2 8C 2.6GHz CPUs with 8 cores each. The examined test problem is based on applying boundary element method for numerical simulation of laminar flow around airfoils [3]. The compression tolerance ϵ , required by the HSS algorithm, is associated with the relative and absolute thresholds ϵ_{rel} and ϵ_{abs} . STRUMPACK allows both to be specified by the user. The compression process stops when one of them is reached.

We present numerical tests for several different settings of ϵ_{rel} (the default value of 10^{-2} as well as 10^{-4} , 10^{-8} and 10^{-12}). For the absolute threshold ϵ_{abs} , the default value of 10^{-8} is used.

The times of the sequential tests are presented on Fig. 1 (left). The asymptotic behavior of $O(n^3)$ is clearly seen for the MKL implementation of the Gaussian elimination. We see also the nearly optimal complexity of the STRUMPACK of the HSS algorithm, as well as the impact of increasing the rank r with the decrease of ϵ_{rel} . STRUMPACK significantly outperforms MKL for all settings of ϵ_{rel} .

The best parallel times are observed when 16 threads are used. The obtained results are plotted on Fig. 1 (right). Not surprisingly, the parallel speedup of Gaussian solver from MKL is better then the STRUMPACK ones. The second important

conclusion concerns the overall performance. For larger n , the observed advantage of STRUMPACK is restricted to the case of relatively larger $\epsilon_{rel} = 10^{-2}, 10^{-4}$. It should be noted that the test with $\epsilon = 10^{-12}$ run slightly faster than the experiments with $\epsilon = 10^{-8}$. This is probably due to achieving more efficient hierarchy in the compressed matrix.

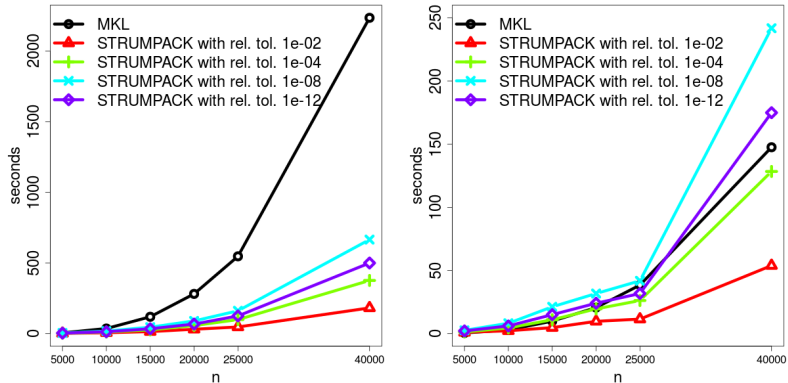


Figure 1: Performance of STRUMPACK and MKL: sequential test (left), and parallel scalability tests using 16 cores (right)

An important question when working with an approximate compression (factorization) method, like HSS, is how accurate it is. In order to examine the used threshold settings of the STRUMPACK package we consider the accuracy in the normalized ℓ_2 norm where the MKL solution is taken as reference/exact.

$$\|x^{Gauss} - x^{HSS}\|_2 = \frac{1}{n} \sqrt{\sum_{i=1}^n (x_i^{Gauss} - x_i^{HSS})^2}$$

In Table 1 we show the calculated ℓ_2 norms varying the problem sizes and the relative tolerances tested. Lowering the threshold improves the accuracy. The results with the lowest relative threshold $\epsilon_{rel} = 10^{-12}$ show the best accuracy and are slightly faster than the tests with the next lower threshold. It should be noted that the accuracy of the results degrades as n grows. With the default threshold of $\epsilon_{rel} = 10^{-2}$, and for sizes $n > 15000$ (the NaN values), the compressed matrix is singular and STRUMPACK cannot solve the system.

4 Concluding remarks

Performance analysis of the STRUMPACK package for solving dense systems of linear algebraic equations arising from the use of Boundary Element Method is presented. A Hierarchical Semi-Separable(HSS) based method is tested on Intel Xeon E5-2650v2 8C 2.6GHz CPUs and the performance and accuracy is studied. The HSS based

n	ϵ_{rel}			
	10^{-2}	10^{-4}	10^{-8}	10^{-12}
5000	1.11e+08	7.10e-03	0.00079	8.22e-06
10000	1.31e+23	1.23e+00	0.0055	1.69e-05
15000	NaN	1.21e+00	0.0050	7.60e-05
20000	NaN	3.69e+05	0.023	8.73e-05
25000	NaN	1.57e+03	0.018	2.28e-04
40000	NaN	1.30e+07	0.001	1.96e-04

Table 1: Normalized ℓ_2 norm of the error

method works significantly faster in the sequential mode but it doesn't scale as well as the direct method from MKL. The accuracy of the method is sensitive with respect to threshold parameters which must be fine tuned for a given size of the problem in order to achieve acceptable results.

References

- [1] P. G. Martinsson. A fast randomized algorithm for computing a hierarchically semiseparable representation of a matrix. *SIAM Journal on Matrix Analysis and Applications*, 32(4):1251–1274, 2011.
- [2] François-Henry Rouet, Xiaoye S. Li, Pieter Ghysels, and Artem Napov. A distributed-memory package for dense hierarchically semi-separable matrix computations using randomization. *ACM Trans. Math. Softw.*, 42(4):27:1–27:35, June 2016.
- [3] Dimitar Slavchev and Svetozar Margenov. Performance analysis of intel xeon phi mics and intel xeon cpus for solving dense systems of linear algebraic equations: Case study of boundary element method for flow around airfoils. In *Proceedings of the 12th Annual Meeting of the Bulgarian Section of SIAM, Studies in Computational Intelligence*. Springer, 2018. Forthcoming.

A Non-Symmetric Model of Disease Propagation

B. Takács, R. Horváth, I. Faragó

One of the oldest problems in the area of applied differential equations is the investigation of the spread of diseases. Since the black death in Europe in the middle ages, people tried to describe and understand the phenomena that dominate these processes.

Since its introduction in the article of Kermack and McKendrick in 1927 ([6]), one of the most popular models has been the SIR model. This model splits the examined population into three disjunctive groups: the susceptible (S), who have yet to contract the disease and become infectious; the infected (I), who can pass on the disease to the previous group; and the recovered (R), who had been infected, but have already recovered and cannot transmit the disease.

The aforementioned model can be written as the following first order system of ordinary differential equations

$$\begin{cases} \frac{dS(t)}{dt} = -aS(t)I(t), \\ \frac{dI(t)}{dt} = aS(t)I(t) - bI(t), \\ \frac{dR(t)}{dt} = bI(t), \end{cases} \quad (1)$$

where $a, b \in \mathbb{R}^+$ are given parameters, and $S(t)$, $I(t)$ and $R(t)$ describe the number of susceptibles, infected and recovered, respectively.

One of the key element of the investigation of the disease is the spatial spread of the infection among the population. This phenomena is characterized by the infection of a single individual, such that in what radius does an infectious person have an effect on the susceptibles. From now on, we will investigate the processes on a two dimensional plane, and later only on a subdomain of it.

Let us expand the previous model (1) in a way that the size of the populations differ in space. For this, we introduce the following system of partial differential equations:

$$\begin{cases} \frac{\partial S(t, x, y)}{\partial t} = -aS(t, x, y)I(t, x, y), \\ \frac{\partial I(t, x, y)}{\partial t} = aS(t, x, y)I(t, x, y) - bI(t, x, y), \\ \frac{\partial R(t, x, y)}{\partial t} = bI(t, x, y), \end{cases} \quad (2)$$

in which $X(x, t)$ denotes the size of population at place x at time t , $X \in \{S, I, R\}$. However, this expansion is not beneficial: the different places behave independently of each other, thus the infection will not spread. In other words, the infection takes place just in a point, i.e. the infected individual at place x only passes on the disease to

those susceptibles who are exactly at x . Instead of this, it is more natural to suppose that the infected individual has an influence on susceptibles in a certain radius around itself, in a way that it infects healthy individuals further from the infected one less likely.

Let the nonnegative function $F(x, y, r, \theta)$ be defined as

$$F(x, y, r, \theta) = \begin{cases} f_1(r)f_2(\theta), & \text{if } (r, \theta) \in B_\delta((x, y)) \\ 0 & \text{otherwise.} \end{cases}$$

where $B_\delta((x, y))$ denotes the δ radius ball with center at (x, y) . This expression measures the influence of the infected at position (x, y) on the susceptibles at (r, θ) , where r and θ are polar-coordinates, i.e. r denotes the radius and θ the angle. Note that in this case we assume a homogeneous domain in which the propagation of the disease does not depend on the place of the infectious individual, i.e. the spread of the disease is the same at every point of the domain (there are no regions in which the disease would spread faster or slower). It is a natural presumption that f_1 decreases as r increases, and f_2 can be either constant in θ , or a periodic function in a way that $f_2(\theta) : [0, 2\pi] \rightarrow \mathbb{R}$, $f_2(0) = f_2(2\pi)$. The case of constant f_2 was widely studied in the articles of I. Faragó and R. Horváth ([3] and [5]). In this talk we consider the nonconstant case, which can be a model of a plant disease propagated by a constant wind on the domain, for example.

Now we would like to expand the original theory with a spatial dependence described by the function F . This way, we get the following equation for the susceptibles:

$$\frac{\partial S(t, x, y)}{\partial t} = - \int_0^\infty \int_0^{2\pi} F(x, y, r, \theta) I(t, x'(r, \theta), y'(r, \theta)) r \, d\theta dr \cdot S(t, x, y), \quad (3)$$

in which we used the notations $x'(r, \theta) = x + r \cos(\theta)$ and $y'(r, \theta) = y + r \sin(\theta)$, and r is the Jacobian determinant.

Note that outside of the previously described $B_\delta((x, y))$ ball the value of this integral is zero. Hence, with the definition of the function F , the previous term takes the form

$$\frac{\partial S(t, x, y)}{\partial t} = - \int_0^\delta \int_0^{2\pi} f_1(r)f_2(\theta) I(t, x'(r, \theta), y'(r, \theta)) r \, d\theta dr \cdot S(t, x, y).$$

This way, the *SIR* model with spatial dependence can be written as the following system of integro-differential equations:

$$\begin{cases} \frac{\partial S(t, x, y)}{\partial t} = - \int_0^\delta \int_0^{2\pi} f_1(r)f_2(\theta) I(t, x'(r, \theta), y'(r, \theta)) r \, d\theta dr \cdot S(t, x, y) \\ \frac{\partial I(t, x, y)}{\partial t} = \int_0^\delta \int_0^{2\pi} f_1(r)f_2(\theta) I(t, x'(r, \theta), y'(r, \theta)) r \, d\theta dr \cdot S(t, x, y) - bI(t, x, y) \\ \frac{\partial R(t, x, y)}{\partial t} = bI(t, x, y) \end{cases} \quad (4)$$

These cannot be solved analytically, so we present here a few numerical methods. It is important to use adequate numerical schemes, which preserve the qualitative properties of the continuous model. In our case, we will consider the following qualitative properties.

C_1 : the numbers of the individuals in classes S, I and R are nonnegative at every point in our domain,

C_2 : the size of the whole population is constant, i.e.

$$S(t, x, y) + I(t, x, y) + R(t, x, y) = N(x, y),$$

for every $(x, y) \in \Omega$, where Ω is the domain we investigate the spread of the disease in, and $N(x, y)$ is constant in time,

C_3 : the size of the population of S is non-increasing in time at every x ,

C_4 : the size of the population of R is non-decreasing in time at every x .

In the talk we consider the Taylor expansion method used in [3] and [5] and show that in the case of a simple nonconstant $f_2(\theta)$ function these are not beneficial. Because of this, we consider a numerical integration to compute the integrals in (4). From now on we use the following notation:

$$F_I(t, x, y, r_i, \theta_j) = f_1(r_i)f_2(\theta_j)r_iI(t, x + r_i \cos(\theta_j), y + r_i \sin(\theta_j)). \quad (5)$$

Let us consider a two dimensional cubature $\mathcal{Q}(x, y)$ on the $B_\delta(x, y)$ disc with positive coefficients:

$$\int_0^\delta \int_0^{2\pi} F_I(t, x, y, r, \theta) d\theta dr \approx \sum_{(r_i, \theta_j) \in \mathcal{Q}(x, y)} w_{i,j} F_I(t, x, y, r_i, \theta_j) =: T(t, \mathcal{Q}(x, y)), \quad (6)$$

in which $w_{i,j} > 0$. Formally,

$$\mathcal{Q}(x, y) := \{(r_i, \theta_j) : (x + r_i \cos(\theta_j), y + r_i \sin(\theta_j)) \in \text{Int}(B_\delta(x, y))\}.$$

With the approximation (6) we get the following differential equation:

$$\begin{cases} \frac{dS(t, x, y)}{dt} = -S(t, x, y)T(t, \mathcal{Q}(x, y)), \\ \frac{dI(t, x, y)}{dt} = S(t, x, y)T(t, \mathcal{Q}(x, y)) - bI(t, x, y), \\ \frac{dR(t, x, y)}{dt} = bI(t, x, y), \end{cases} \quad (7)$$

Theorem 1. *Properties C_1, C_2, C_3 and C_4 hold without any restrictions.*

For our numerical solutions we will split our domain using a spatial grid, and approximating the continuous solutions by a vector of the values at the gridpoints. We also approximate the derivatives on the left size of (7), and consequently get a numerical scheme. In the talk the discrete version of Theorem 1 will be stated for two different numerical schemes derived from (7). Also, several numerical experiments will also be presented.

References

- [1] COOLS, R., KIM, K.J. , *A survey of known and new cubature formulas for the unit disc*, Korean J. Comput. and Appl. Math. Vol. 7 (2000), No. 3, pp. 477 - 485
- [2] DAVIS, P., POLONSKY, I., *Numerical Interpolation, Differentiation and Integration*, in: M. ABRAMOWITZ, I. STEGUN, *Handbook of Mathematical Functions*, NBS, (1964), pp. 876-925.
- [3] FARAGÓ, I., HORVÁTH, R., *On some qualitatively adequate discrete space-time models of epidemic propagation*, Journal of Computational and Applied Mathematics, 293 (2016) 45-54
- [4] FARAGÓ, I., HORVÁTH, R., *Qualitative properties of the finite difference solution of a space-time epidemic propagation model*, Annales Univ. Sci. Budapest., Sect. Comp. 45 (2016) 157-168
- [5] FARAGÓ, I., HORVÁTH, R., *Qualitative properties of some discrete models of disease propagation*, Journal of Computational and Applied Mathematics (2017)
- [6] W.O. KERMACK, A.G. MCKENDRICK, *A contribution to the mathematical theory of epidemics*, Proc. R. Soc. A: Math. Phys. Eng. Sci. 115 (772) (1927) 235240.

Artificial Neural Networks Time Series Forecasting with Android Live Wallpaper Technology

P. Tomov, I. Zankinski, M. Barova

Abstract

The application of Artificial Neural Networks (ANN) for time series forecasting is quite common in the last few decades [1]. ANN may be adopted in a variety of topologies and training algorithms. Algorithms can be sequential, but they can also be performed in parallel. Computing device types used for training ANN may vary from supercomputers, grid networks and single desktop computers to mobile devices. The focus of this research is training of ANN for Time Series Forecasting (TSF) as background calculations of Android Live Wallpaper technology.

Keywords: artificial neural networks, mobile computing, time series forecasting

1 Introduction

In their nature, common types of ANNs are weighted directed graphs [5]. The process of ANN training aims to find such values for the weights which minimize the total error of the output [6]. When the size of ANN is larger, the amount of weights may reach very high levels. With the increase of the number of weights, the velocity of training of a single ANN on a single computer is being reduced. The elevated amount of calculations poses the application of supercomputers, clusters, grids and donated distributed computing networks as reasonable and adequate [3, 4, 7, 10]. In the last decade the immense spread of mobile devices suggests their surpassing usage over stand alone computers. During most of the operating time, the mobile devices are in idle mode. The possibilities ANN to be trained in distributed computing environment and the use of mobile devices can be efficiently combined in software system for ANN training on mobile operating systems such as Android. In the current research Android Live Wallpaper technology is involved in ANN training for time series forecasting.

2 Technical Solution

The technical solution is divided in two parts. The first one is related to the Android application itself. The second one is related to the ANN data structure representation, information representation and the process of training/forecasting. For the first part Android development capabilities are used while for the second part Encog programming library is involved.

2.1 Android Live Wallpaper

Live Wallpaper is interactive background on the Android home screen, which can even process animation. The Live Wallpaper does not differ much from other Android applications. First step of Live Wallpaper creation is an XML file with application description (manifest). Second component is a Service class inherited from WallpaperService. Inside of the Service class there is an Engine class. The Engine object is responsible for training/forecasting execution and background redrawing. Running LW on the Android operating system requires special permissions. The use of Live Wallpaper feature should be explicitly written into the manifest file in order to prevent wallpaper installation on devices which are not capable of running it.

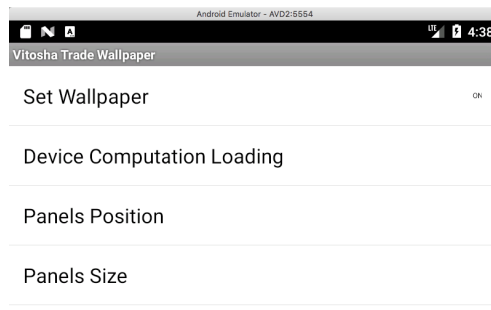


Figure 1: Android Live Wallpaper settings screen.

The set up the wallpaper is performed by sending an Intent to the operating system. SQLite database is used to store financial time series in the mobile device for offline mode operation. On regular intervals, wallpaper service is activated - cycle of training is executed, forecast is retrieved and the visual information is updated. Settings screen is used for visual representation parameters and device loading options (Fig. 1). During operating mode a background image is drawn and ANN training/forecasting information is displayed (Fig. 2).

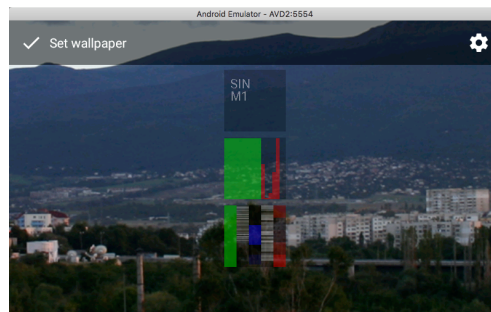


Figure 2: Android Live Wallpaper operational screen.

2.2 Encog Machine Learning Framework

For the process of forecasting, Econg library is used. Multilayer perceptron with 3 layers is chosen. Time series is conditionally divided in two parts (past and future). Data frame for the past (lag) is supplied at the input of the ANN. Data frame for the the future (lead) is expected at the output of the ANN. Time series data are scaled (according neurons activation function) before fed into ANN's input. The output is also scaled with the opposite scaling coefficient used at the input. At each training cycle resilient backpropagation training is executed.

3 Conclusions

As per the results of the current research, the application of donated mobile devices power appears to be much more promising even compared to the donated desktop distributed computing, given that mobile devices are almost always running, which is not the case with the desktop computers. As future studies, donated mobile distributed computing infrastructure can be efficiently used for experiments with different activation functions [11] or permutation algorithms [12]. Other interesting research areas where mobile distributed computing can be applied are barcode readers [2] and computer networks traffic analysis [8, 9].

Acknowledgements

This work was supported by private funding of Velbazhd Software LLC.

References

- [1] Atanasova, T., Barova, M., Balabanov, T., *Use of Neural Models for Analysis of Time Series in Big Data*, Publishing complex of "Vasil Levski" National Military University, ISSN 1314-1937, 193–198, 2016.
- [2] Atanasova, T., Atanasov, J., *Business Processes Traceability in SME by Barcode System*, Proceedings of the International Scientific Conference, UNITECH16, Gabrovo, Bulgaria, ISSN 1313-230X, 207–212, 2016.
- [3] Balabanov, T., Genova, K., *Distributed System for Artificial Neural Networks Training Based on Mobile Devices*, Proceedings of the International Conference Automatics and Informatics, Sofia, Bulgaria, Federation of the Scientific Engineering Unions John Atanasoff Society of Automatics and Informatics, ISSN 1313-1850, 49–52, 2016.
- [4] Balabanov, T., Keremedchiev, D., Goranov, I., *Web Distributed Computing For Evolutionary Training Of Artificial Neural Networks*, International Conference InfoTech, Varna - St. St. Constantine and Elena resort, Bulgaria, Publishing House of Technical University - Sofia, ISSN 1314-1023, 210–216, 2016.

- [5] Balabanov, T., Zankinski, I., Barova, M., *Strategy for Individuals Distribution by Incident Nodes Participation in Star Topology of Distributed Evolutionary Algorithms*, Cybernetics and Information Technologies, Institute of Information and Communication Technologies - BAS, vol. 16, no. 1, ISSN 1311-9702, 80–88, 2016.
- [6] Balabanov, T., Zankinski, I., Dobrinkova, N., *Time Series Prediction by Artificial Neural Networks and Differential Evolution in Distributed Environment*. Proceedings of the International Conference on Large-Scale Scientific Computing, Sozopol, Bulgaria, Lecture Notes in Computer Science, Springer, vol. 7116, no. 1, ISBN 978-3-642-29842-4, 198205, 2011.
- [7] Keremedchiev, D., Barova, M., Tomov, P., *Mobile Application as Distributed Computing System for Artificial Neural Networks Training Used in Perfect Information Games*, Proceedings of the International Scientific Conference, UNITECH16, Gabrovo, Bulgaria, ISSN 1313-230X, 389–393, 2016.
- [8] Tashev, T., Marinov, M., Monov, V., Tasheva, R., *Modeling of the MiMa-algorithm for crossbar switch by means of Generalized Nets*, Proceedings of the 2016 IEEE 8th International Conference on Intelligent Systems (IS), Sofia, Bulgaria, ISBN 978-1-5090-1354-8, 593–598, 2016.
- [9] Tashev, T., Monov, V., *Modeling of the hotspot load traffic for crossbar switch node by means of generalized nets*, Proceedings of the 6-th International IEEE Conference Intelligent Systems IS'12, Sofia, Bulgaria, vol. 2, 187–191, 2012.
- [10] Tomov, P., Monov, V., *Artificial Neural Networks and Differential Evolution Used for Time Series Forecasting in Distributed Environment*, Proceedings of the International Conference Automatics and Informatics, Sofia, Bulgaria, ISSN 1313-1850, 129–132, 2016.
- [11] Zankinski, I., Tomov, P., Balabanov, T., *Alternative Activation Function Derivative in Artificial Neural Networks*, 25th Symposium with International Participation - Control of Energy, Industrial and Ecological Systems, Bankya, Bulgaria, John Atanasoff Union of Automation and Informatics, ISSN 1313-2237, 79–81, 2017.
- [12] Zankinski, I., Stoilov, T., *Effect of the Neuron Permutation Problem on Training Artificial Neural Networks with Genetic Algorithms in Distributed Computing*, Proceedings of the 24th International Symposium Management of Energy, Industrial and Environmental Systems, ISSN 1313-2237, Bankya, Bulgaria, 53–55, 2016.

Pareto Optimal Solutions of Noise Statistics for Kalman Filtering Applied to State Estimation of Gas Dynamics

F. E. Uilhoorn

We present an optimization framework that seeks the Pareto optimal solutions of the model noise statistics for the extended Kalman filter (EKF) applied to state estimation of gas flow dynamics in pipelines. Recursive Bayesian estimators, like the EKF (see Algorithm 1) enable us to combine noisy measurements with a simulator that solves an inexact flow model [2, 6, 11, 10, 12]. Mathematically, we read

$$\xi_k = f(\xi_{k-1}) + v_{k-1}, \quad (1)$$

$$y_k = h(\xi_k) + n_k, \quad (2)$$

where $v_{k-1} \sim \mathcal{N}(0, Q_{k-1})$ and $n_k \sim \mathcal{N}(0, R_k)$ with model and measurement noise covariance matrices Q_{k-1} and R_k , respectively. Since, we assume an isothermal flow field, the state vector ξ_k represents pressure p and mass flow rate \dot{m} . The functions f and h refer to the flow and measurement model, respectively.

ALGORITHM 1: EXTENDED KALMAN FILTER.

```

1: procedure EKF_ALGORITHM( $y_{1:n_k}$ )
2:    $\hat{\xi}_{0|0} \leftarrow \xi_0$  and  $P_{0|0} \leftarrow P_0$ 
3:   for  $k = 1 : n_k$  do
4:      $\hat{\xi}_{k|k-1} \leftarrow f_{k-1}(\hat{\xi}_{k-1|k-1})$ 
5:      $F_{k-1} \leftarrow [\nabla_{\xi_{k-1}} f_{k-1}^\top(\xi_{k-1})]^\top$ 
6:      $P_{k|k-1} \leftarrow F_{k-1} P_{k-1|k-1} F_{k-1}^\top + Q_{k-1}$ 
7:      $H_k \leftarrow [\nabla_{\xi_k} h_k^\top(\xi_k)]^\top$ 
8:      $K_k \leftarrow P_{k|k-1} H_k^\top (R_k + H_k P_{k|k-1} H_k^\top)^{-1}$ 
9:      $\hat{\xi}_{k|k} = \hat{\xi}_{k|k-1} + K_k (y_k - h_k(\hat{\xi}_{k|k-1}))$ 
10:     $P_{k|k} = (I - K_k H_k) P_{k|k-1} (I - K_k H_k)^\top + K_k R_k K_k^\top$ 
11:   end for
12: end procedure

```

It is of key importance to identify the optimal process noise statistics because it can have a significant impact on the filtering accuracy. Improper specification of the noise variances can cause filter divergence. A priori knowledge about these noise statistics is in practice unavailable. In fact, the elements of Q , also called filter tuning parameters, are often obtained via an ad-hoc process based on trial-and-error [7].

Based on the notion of Pareto optimality, our aim is to find a trade-off between incommensurable objective functions that are related to estimation errors of p and \dot{m} . Considering a multiobjective problem we can write

$$\min_{\psi \in \mathcal{X}} F(\psi) = (f_1(\psi), f_2(\psi), \dots, f_q(\psi)), \quad (3)$$

where $F : \mathbb{R}^n \rightarrow \mathbb{R}^q$ with n and q as the number of variables and objective functions, respectively. Let $\mathcal{X} = \{\psi \in \mathbb{R}^n : a \leq \psi \leq b\}$ where a and b are the bound vectors. For the isothermal flow model (5) we formulated two single objective functions. Thus, $q = 2$ and $\psi \in \{q_{p,11}, q_{p,22}, \dots, q_{p,nn}, q_{\dot{m},11}, q_{\dot{m},22}, \dots, q_{\dot{m},nn}\}$ where q_p and $q_{\dot{m}}$ denote the elements of covariance matrix Q and n is the number of measurement points. The performance index to be minimized is based on the mean spatial and temporal root mean square error. For p and in the same manner for \dot{m} , it is defined as

$$f_1(q_{p,11}, q_{p,22}, \dots, q_{p,nn}) = \left(\frac{\|X_p - \hat{X}_p\|_F}{\sqrt{n_x n_k}} \right), \quad (4)$$

where X is the true and \hat{X} is the estimated matrix and within the discretized domain n_x and n_k are the number of nodes and time steps, respectively. The measurement noise statistics represented by R in the EKF were assumed to be known and in practice often obtained from sensor testing and calibration. The optimization problem was solved using the Bi-Objective Mesh Adaptive Direct Search (BiMADS) [1] because it has the convenience, that no information is needed about the gradient or even higher derivative to search for the optimal solution (see Algorithm 2). The isothermal gas flow model is described as follows

$$\frac{\partial u}{\partial t} + \frac{\partial f(u)}{\partial x} = s(u), \quad \forall x \in [0, L], t \in [0, t_f], \quad (5)$$

with $u = [p \ \dot{m}]^\top$, $f(u) = [a_s^2 \dot{m} A^{-1} \ Ap]^\top$ and $s(u) = [0 \ -fa_s^2 \dot{m} |\dot{m}| (2dAp)^{-1}]^\top$. The flow model was approximated with a semidiscrete finite volume scheme using Roe's flux limiter. The approximation of $u(x, t)$ can be done as follows

$$\frac{d}{dt} U_i(t) + \frac{1}{\Delta x_i} (\mathcal{F}_{i+1/2} - \mathcal{F}_{i-1/2}) = s_i, \quad (6)$$

where $\mathcal{F}_{i+1/2} = \mathcal{F}(U_{i+1/2}^-, U_{i+1/2}^+)$ is the monotone numerical flux that approximates $f(u(x_{i+1/2}, t))$. The approximations $U_{i+1/2}^-, U_{i+1/2}^+$ of the point value $u(x_{i+1/2}, t)$ are obtained via a reconstruction process using Roe's superbee flux limiter [8]. This to minimize the presence of numerical dissipation. The resulting system of ODEs was integrated with the explicit Runge–Kutta scheme. The Jacobian in EKF was approximated by a finite difference scheme. The elements of the Jacobian are $((\partial dx_i / dt) / \partial x_j) = \partial f_{k-1, [i]} / \partial x_{k-1, [j]}$ where dx_i / dt refers to the i^{th} ODE with x representing p and \dot{m} .

Numerical experiments were conducted with boundary condition $\dot{m}(L, t)$ containing gradual changes and a hydraulic shock. First, we assumed that the measurement noise for each instrument is constant. For this idealized situation, we compared BiMADS with the normalized WS method [3, 9] and widely used Non-Dominated Sorting Genetic Algorithm (NSGA-II) [5]. For the WS method, the ideal vector z^* was found by $z_i^* = \min_{\psi \in \mathcal{X}} f_i(\psi)$ and the nadir point by $z_i^{\text{nad}} = \max_{\psi \in \mathcal{P}} f_i(\psi)$. This method resulted in duplicated solutions. The evolution algorithm NSGA-II was

ALGORITHM 2: BiMADS.

```

1: procedure BiMADS( $f_1, f_2, \mathcal{X}, \psi_0$ )
2:   Solve  $\min_{\psi \in \mathcal{X}} f_j(\psi)$ ,  $j \in \{1, 2\}$  with MADS.
3:   Sort nondominated points  $\mathcal{X}_{\mathcal{L}}$  in ascending order of  $f_1$  and descending order of  $f_2$ .
4:    $w(\psi) \leftarrow 0 \forall \psi \in \mathcal{X}$  and  $\delta > 0$ .
5:   for  $k = 0, 1, 2, \dots$  do
6:     if  $L_k = 1$  then
7:        $\psi^{\hat{j}} \leftarrow \psi^1$ ,  $\delta^{\hat{j}} \leftarrow \delta / (w(\psi^{\hat{j}}) + 1)$ 
8:       solve  $\min_{\psi \in \mathcal{X}} f_j(\psi)$ ,  $j \in \{1, 2\}$  with MADS from  $\psi^{\hat{j}}$ .
9:     else if  $L_k = 2$  then
10:       $\psi^{\hat{j}} \leftarrow \psi^2$ ,  $r \leftarrow f_1((\psi^2), f_2(\psi^1))$  and  $\frac{\|F(\psi^2) - F(\psi^1)\|^2}{w(\psi^2) + 1}$ .
11:     else if  $L_k > 2$  then
12:       $\hat{j} \in \operatorname{argmin} \delta^j \leftarrow \frac{\|F(\psi^j) - F(\psi^{j-1})\|^2 + \|F(\psi^j) - F(\psi^{j+1})\|^2}{w(\psi^j) + 1}$ .
13:       $r \leftarrow f_1((\psi^{\hat{j}+1}), f_2(\psi^{\hat{j}-1}))$ .
14:     end if
15:     Solve single-objective function  $\min_{\psi \in \mathcal{X}} \Psi_r(\psi)$  from  $\psi^{\hat{j}}$  with MADS.
16:     Update the set of nondominated points  $\mathcal{X}_{\mathcal{L}}$  by adding new ones.
17:     Remove dominated points and order list of points.
18:      $w(\psi^{\hat{j}}) \leftarrow w(\psi^{\hat{j}}) + 1$  for  $\psi \in \mathcal{X}_{\mathcal{L}}$ .
19:   end for
20: end procedure

```

not only significantly slower but also less efficient compared to BiMADS. Hence, the succeeding computations were only done with the latter algorithm. Based on the concept of normal boundary intersection [4], the knee point was calculated to obtain the optimal values of q required for state estimation. In the more realistic scenario, simulations were conducted whereas each instrument along the pipeline had different noise parameters. Till now, the model noise was assumed identical, therefore in the last step, we seek to optimize q at each point for p and \dot{m} . The simulations were conducted for gradual changes in $\dot{m}(L, t)$, therefore in the final experiment, we imposed a hydraulic shock and repeated the preceding simulations. Results showed that BiMADS is suitable for designing the EKF algorithm for estimating the gas flow dynamics. Significant higher computation times were recorded when the number of tuning variables increased.

References

- [1] C. Audet, G. Savard, and W. Zghal. Multiobjective optimization through a series of single-objective formulations. *SIAM J Optim*, 19(1):188–210, 2008b.
- [2] H.A. Behrooz and R. B. Boozarjomehry. Modeling and state estimation for gas transmission networks. *Journal of Natural Gas Science and Engineering*, 22:551–570, 2015.

- [3] J. L. Cohon. *Multiobjective programming and planning*. New York: Academic Press, 1978.
- [4] I. Das. On characterizing the knee of the Pareto curve based on normal-boundary intersection. *Structural and Multidisciplinary Optimization*, 18(2):107–115, 1999.
- [5] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan. A fast and elitist multiobjective genetic algorithm: NSGA-II. *IEEE Transactions on Evolutionary Computation*, 6(2):182–197, 2002.
- [6] I. Durgut and K. Leblebicioğlu. Kalman-Filter Observer Design around Optimal Control Policy for Gas Pipelines. *International Journal for Numerical Methods in Fluids*, 24(2):233–245, 1997.
- [7] T. D. Powell. Automated Tuning of an Extended Kalman Filter Using the Downhill Simplex Algorithm. *Journal of Guidance, Control, and Dynamics*, 25(5):901–908, 2002.
- [8] P.L. Roe. Characteristic-based schemes for the Euler equations. *Annual Review of Fluid Mechanics*, 18(1):337–365, 1986.
- [9] S. Shan and G. G. Wang. An efficient pareto set identification approach for multiobjective optimization on black-box functions. *Journal of Mechanical Design*, 127(5):866, 2005.
- [10] F. E. Uilhoorn. Estimating rapid flow transients using extended Kalman filter. *Silesian J. Pure Appl. Math.*, 6(1):97–110, 2016.
- [11] F. E. Uilhoorn. Comparison of Bayesian estimation methods for modeling flow transients in gas pipelines. *Journal of Natural Gas Science & Engineering*, 38:159–170, 2017.
- [12] F.L.V. Vianna, H.R.B. Orlande, and G.S. Dulikravich. Estimation of the temperature field in pipelines by using the Kalman filter. In *2nd International Congress of Serbian Society of Mechanics: Serbia*, 2009.

Implementation of the Three-times Repeated Richardson Extrapolation together with Explicit Runge-Kutta Methods

Z. Zlatev, I. Dimov, I. Farago, K. Georgiev, A. Havasi

1 Introduction of the Three-times Repeated Richardson Extrapolation

The Three-times Repeated Richardson Extrapolation can successfully be combined with Explicit Runge-Kutta Methods (ERKMs) and used in the numerical treatment of non-linear systems of ordinary differential equations (ODEs). These combinations are new numerical methods for solving systems of ODEs, which have to be studied carefully, which must be done and is very essential with regard to their stability properties. The computational cost per step of the new numerical methods is higher than the computational cost per step of the underlying ERKMs.

Consider non-linear systems of first order ordinary differential equations (ODEs) defined as follows:

$$\frac{dy}{dt} = f(t, y), \quad t \in [a, b], \quad a < b, \quad y \in \mathbf{R}^s, \quad f \in D \subset \mathbf{R}^s \times \mathbf{R}^s, \quad y(a) = \eta \quad (1)$$

and assume that these systems are to be solved on the following set of equidistant grid-points:

$$t_0 = a, \quad t_n = t_{n-1} + h, \quad (n = 1, 2, \dots, N), \quad t_N = b, \quad h = \frac{b-a}{N}, \quad (2)$$

by applying an *arbitrary one-step numerical method*, the order of accuracy of which is p . The fact that one-step numerical methods are used means that *only* the approximation $y_{n-1} \approx y(t_{n-1})$ is used in the calculation of the next approximation $y_n \approx y(t_n)$ for any $n \in 1, 2, \dots, N$. Assume furthermore that the computations at the points t_1, t_2, \dots, t_{n-1} of the grid (2) are completed and that the calculations at point t_n have to be carried out. Much more details about the one-step methods can be found for example in [2]. If the Three-times Repeated Richardson Extrapolation is used with an arbitrary one-step method, then the approximation y_n can be computed by using the previous approximation y_{n-1} and by performing successively the following six calculation processes:

1. Apply the selected ERKM to compute an approximation $z_n^{[1]}$ of the solution of (1) at the point $t = t_n$ by performing one step with a stepsize h .
2. Apply the selected ERKM to compute an approximation $z_n^{[2]}$ of the solution of (1) at the point $t = t_n$ by performing two steps with a stepsize $h/2$.

3. Apply the selected ERKM to compute an approximation $z_n^{[3]}$ of the solution of (1) at the point $t = t_n$ by performing four steps with a stepsize $h/4$.
4. Apply the selected ERKM to compute an approximation $z_n^{[4]}$ of the solution of (1) at the point $t = t_n$ by performing eight steps with a stepsize $h/8$.
5. Apply the selected ERKM to compute an approximation $z_n^{[5]}$ of the solution of (1) at the point by performing sixteen steps with a stepsize $h/16$.
6. Compute an approximation y_n of the solution of (1) at the point $t = t_n$ by using the approximations $z_n^{[k]}$, $k = 1, \dots, 5$ obtained in the previous five calculation processes.

It should be mentioned here that the Classical Richardson Extrapolation (based essentially on the first two calculation processes) was first introduced in [3]; see more details in [6, 7], the Repeated Richardson Extrapolation (based essentially on the first three processes) was studied in [4], while the Two-times Repeated Richardson Extrapolation (based on the first four processes) is described in [5].

2 Accuracy of the Three-times Repeated Richardson Extrapolation

The order of accuracy of the combined methods is much higher: if the order of accuracy of the underlying EPKM is p then the order of accuracy of its combination with a Three-times Repeated Richardson Extrapolation is at least $p + 4$ (under the assumption that the right-hand-side function of the system of ODEs is sufficiently many times continuously differentiable).

The following theorem is proven:

Theorem *Consider the numerical solution of the system of ODEs (1) that is obtained by using an arbitrary one-step method and assume that the chosen method is of order of accuracy p . Then the order of accuracy of the combination of the Three-times Repeated Richardson Extrapolation and the chosen one-step method is at least $p+4$ when the function $f(t, y)$ from the right-hand side of (1) is at least $p+4$ times continuously differentiable.* The result proved above shows that the order of accuracy can be in-

creased very significantly when the *Three-times Repeated Richardson Extrapolation* is used, but it is necessary to pay some price (to carry out thirty-one steps instead of only one) for computing the accurate approximation y_n . In our full length paper we demonstrate, by using two numerical examples, that the high accuracy of the Three-times Repeated Richardson Extrapolation is sometimes allowing us to increase the stepsize and, by solving the problem (1) with a sufficiently large stepsize, to achieve both the required accuracy and a very good compensation for the need to use much more computations in the performance of the six calculation processes described in Section 1.

3 ERKMs combined with the Three-times Repeated Richardson Extrapolation

The two advantages, higher accuracy and better stability, are often giving a very good compensation for the increased computational cost per step, because the same degree of accuracy can be achieved by applying a large stepsize which leads to a considerable reduction of the number of steps when the Three-times Repeated Richardson Extrapolation is used in combination with ERKMs.

In this section we introduce: (a) the Explicit Runge-Kutta Methods (the ERKMs), (b) their stability polynomials together with the stability polynomials of the Three-times Repeated Richardson Extrapolation and (c) their absolute stability regions again together with the absolute stability regions of the Three-times Repeated Richardson Extrapolation (for comparison the absolute stability regions of the Classical Richardson Extrapolation, the Repeated Richardson Extrapolation and the Three-times Repeated Richardson Extrapolation will also be given). The class of the ERKMs is a sub-class of the one-step numerical methods.

4 Selecting particular Explicit Runge-Kutta Methods

Some particular numerical methods satisfying the conditions $p = m$ and $m = 1, 2, 3, 4$ are needed for the considered numerical experiments. Such methods is presented in this section. If $p = m = 1$, then only one Explicit Runge-Kutta Method exists, the Forward Euler Formula. For each pair p, m with $p = m$ and $m = 2, 3, 4$ there exists a large class of Explicit Runge-Kutta Methods (depending on one or two free parameters for $m = 2$ and $m = 3, 4$ respectively; see again [2]). All methods from such a class have the same absolute stability region. Furthermore, each of these methods can be used in combination with any version of the Richardson Extrapolation and, if the version is fixed, then all such combinations have the same absolute stability region. We shall choose particular methods for each pair p, m with $p = m$ and $m = 2, 3, 4$. The selected methods which are used to run the numerical examples are:

- (a) One-stage first-order Explicit Runge-Kutta Method
- (b) Two-stages second-order Explicit Runge-Kutta Method
- (c) Three-stages third-order Explicit Runge-Kutta Method
- (d) Four-stages fourth-order Explicit Runge-Kutta Method

5 General conclusions

The Three-times Repeated Richardson Extrapolation is giving *very accurate results*. Its order of accuracy is at least $p + 4$ when the order of accuracy of the underlying

method is only p . This property often allows us to use large stepsizes and to reduce considerably the computational cost, achieving at the same time much higher accuracy than that obtained by using the ERKMs (and also much higher accuracy than that achieved when the Classical Richardson Extrapolation, the Repeated Richardson Extrapolation and the Two-times Repeated Richardson Extrapolation are used).

It was demonstrated that (a) the Three-times Repeated Richardson Extrapolation have *better stability properties* than those of the ERKMs when $p = m$ and $m = 1, \dots, 4$ and (b) this device can sometimes produce stable results in cases where the underlying ERKM is unstable. However, it should be emphasized here that only the better stability properties of the Three-times Repeated Richardson Extrapolation will in general not result in more efficient computational processes. The fact that high accuracy can be obtained by specifying larger stepsizes during the computations should additionally be exploited in the efforts for achieving greater efficiency.

It was assumed in this paper that a *constant stepsize* is used during the calculations. Strictly speaking, such an assumption is not needed (not always, at least), but the exposition of the results was facilitated considerably by making it.

The use of the Classical Richardson Extrapolation, the Repeated Richardson Extrapolation and the Two-times Repeated Richardson Extrapolation together with Explicit Runge-Kutta Methods is only briefly mentioned in this paper, but we have presented the absolute stability regions of the methods based on the use of the Classical Richardson Extrapolation, the Repeated Richardson Extrapolation and the Two-times Repeated Richardson Extrapolation as well as numerical results obtained by applying these three numerical methods. Much more details about the application of the Classical Richardson Extrapolation, the Repeated Richardson Extrapolation and the Two-times Repeated Richardson Extrapolation together with Explicit Runge-Kutta Methods can be found in the second chapter of the monograph [6], see also [4, 5].

References

- [1] G. Dahlquist: *A special stability problem for linear multistep methods*, BIT, Vol. 3 (1963), pp. 27-43.
- [2] J. D. Lambert: *Numerical Methods for Ordinary Differential Equations: The Initial Values Problem*, Wiley, New York, 1991.
- [3] L. F. Richardson: *The Deferred Approach to the Limit, ISingle Lattice*, Philosophical Transactions of the Royal Society of London, Series A, Vol. 226 (1927), pp. 299-349.
- [4] Z. Zlatev, I. Dimov, I. Farag, K. Georgiev and A. Havasi, *Stability properties of the Repeated Richardson Extrapolation combined with some Explicit Runge-Kutta Methods*, Talk presented at the SIAM Conference in Sofia, December 22 2017.
- [5] Z. Zlatev, I. Dimov, I. Farago, K. Georgiev and A. Havasi, *Stability properties of the Two-times Repeated Richardson Extrapolation combined with some*

Explicit Runge-Kutta Methods, Talk presented at the conference in Losentz, June 2018.

- [6] Z. Zlatev, I. Dimov, I. Farago and A. Havasi: *Richardson Extrapolation: Practical Aspects and Applications*, De Gruyter, Berlin, 2017.
- [7] Z. Zlatev, K. Georgiev and I. Dimov, *Studying absolute stability properties of the Richardson Extrapolation combined with Explicit Runge-Kutta Methods*, Computers and Mathematics with Applications, Vol. 67, No. 12 (2014), pp. 2294-2307.

An adaptive Newton method for solving nonlinear partial differential equations

O. Axelsson, S. Sysala

The convergence of the full step Newton iteration method requires that already the initial approximation is sufficiently close to the exact solution. The normal procedure to cope with that is to use a damped step version of the method, which however complicates the method as it needs many iteration steps and can make it converge very slow. The damped step version can be seen as the lowest order time-stepping method to solve the corresponding evolution equation to find the stationary solution. In this paper we consider higher order, namely a second order and a third order time stepping method and show that then, in particular the third order method has favourable properties, i.e. much faster convergence and can converge for an initial approximation in a much larger ball around the exact solution. We also present an efficient adaptive mesh refinement procedure to gradually approach the solution for a sufficiently fine mesh. The methods are illustrated numerically for an elastoplastic problem.

The normal Newton method to solve a nonlinear equation $F(x) = 0$ has the form

$$F'(x^k)(x^{k+1} - x^k) = -\tau_k F(x^k), \quad k = 0, 1, \dots$$

where x^0 is given and $\{\tau_k\}$ is a set of time-step parameters. The method converges in general only if the set $\{\tau_k\}$ is sufficiently small. If convergence holds for $\tau_k = 1$, the method converges superlinearly, normally with a quadratic rate. The method can be seen as a time-stepping method for the evolution equation

$$\frac{d}{dt}F(x(t)) = -F(x(t)), \quad t > 0, \quad x(0) = x_0.$$

Since $F(x(t)) = e^{-t}F(x(0))$, i.e. the solution converges exponentially to the stationary solution x_∞ of $F(x_\infty) = 0$, where $x(t) \rightarrow x_\infty$, $t \rightarrow \infty$ and $x_\infty = \hat{x}$, the solution of $F(\hat{x}) = 0$, it indicates that one can get convergence even if one uses full time-steps, $\tau_k = 1$, at least if one uses higher order time-integration methods. We analyze this for a second and third order method and show that the convergence is not only very rapid if it converges, namely with a quartic and sixteenth order rates. Furthermore, the initial solution can be located in a larger ball around the exact solution than for the standard first order method. Clearly one can construct still even faster convergent methods in this way.

The second part of the paper shows how one can approach the solution of a nonlinear partial differential equation, normally or elliptic type, by a sequence of steadily refined meshes. Thereby the idea is to solve first the nonlinear problem on a coarse mesh, which is normally relatively cheap and then interpolate the solution on a fine mesh, where it can suffice to use just one linearized solution to get an approximate solution of the same order of accuracy as the fine mesh discretized solution has. This idea was

first presented in the beginning of the 1990'th by J. Xu [1, 2] and O. Axelsson [3] for basic types of nonlinear elliptic problems.

It has later been extended to more general problems, see for instance [4]. However, as shown in [5, 6], the method is a bit too simple in the sence that the above favourable property holds only for sufficiently smooth problems. For problems with local singularities one must use more refinements then just one coarse and a single fine mesh. Such mesh refinements can be done adaptively by use of local refinements where the residuals of the given problem is relatively larger than the average residual. In [6] this was illustrated by first rewriting higher order operators. The Lipschitz condition was not of a relative type.

In this paper we show how an adaptive method can be constructed based on a relative Lipschitz condition. We show also how one can use a sequence of only locally refined meshes to gradually approach an approximate solution with sufficiently small residuals. Thereby in a 2D problem, the mesh refinements is based on a sequence of isosceles triangular which are successively and only locally refined by use of the longest bisection refinement method.

After the construction of a new refined mesh one solves the linearized Newton equation, using the interpolate of the previous approximation as initial value. Since there are few additional points in the mesh, very few, typically just two Newton steps need to be computed on the refined mesh.

Eventually, when all residuals are sufficiently small the whole process can stop.

For an earlier use of adaptive meshes, see [7, 8]. For the use of relative Lipschitz conditions and an alternative continuation Newton method, see [9].

References

- [1] J. Xu. A novel two-grid method for semilinear elliptic equations. *SIAM J. Sci. Comput.* 15(1994) 231–237.
- [2] J. Xu. Two-grid discretization techniques for linear and nonlinear PDEs. *SIAM J. Numer. Anal.* 33(1996) 1759–1777.
- [3] O. Axelsson. On mesh independence and Newton-type methods. *Appl. Math.* 38(1993) 249–265.
- [4] O. Axelsson, W. Layton. A two-level discretization of nonlinear boundary value problems. *SIAM J. Numer. Anal.* 33(1996) 2359–2374.
- [5] I.E. Kaporin, O. Axelsson. On a class of nonlinear equation solvers based on the residual norm reduction over a sequence of affine subspaces. *SIAM J. Sci. Comput.* 16(1995) 228–249.
- [6] O. Axelsson, I.E. Kaporin. Minimum residual adaptive multilevel finite element procedure for the solution of nonlinear stationary problems. *SIAM J. Numer. Anal.* 35(1998) 1213–1229.

- [7] W.F. Mitchell. Optimal multilevel iterative methods for adaptive grids. *SIAM J. Sci. Comput.* 13(1992) 146–167.
- [8] K. Eriksson, C. Johnson. Adaptive streamline diffusion finite element methods for stationary convection-diffusion problems. *Math. Comp.* 60(1993) 167–188.
- [9] O. Axelsson, S. Sysala. Continuation Newton methods. *Computers and Mathematics with Applications* 70(2015) 2621–2637.

Part B

List of participants

Peter Arbenz

ETH Zurich
Computer Science Department
Universitätstrasse 6 CAB F51.1
8092 Zurich, Switzerland
arbenz@inf.ethz.ch

Regina Arbenz

ETH Zurich
Computer Science Department
Universitätstrasse 6 CAB F51.1
8092 Zurich, Switzerland

Owe Axelsson

Institute of Geonics CAS
Studentska 1768
70800 Ostrava, Czech Republic
owea@it.uu.se

Todor Balabanov

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bontchev Str., bl. 2
1113, Sofia, Bulgaria
todor.balabanov@gmail.com

Gergana Bencheva

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bontchev Str., bl. 25A
1113, Sofia, Bulgaria
gery@parallel.bas.bg

Radim Blaheta

Institute of Geonics CAS
Studentska 1768
70800 Ostrava, Czech Republic
blaheta@ugn.cas.cz

Aycil Cesmelioglu

Oakland University
146 Library Drive 368 MSC
48309 Rochester, United States
cesmelio@oakland.edu

Ivan Dimov

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 25A
1113 Sofia, Bulgaria
ivdimov@bas.bg

István Faragó

Eötvös Loránd University &
MTA-ELTE NumNet Research Group
Pázmány P. sétány 1/c
1117 Budapest, Hungary
faragois@cs.elte.hu

Georgi Gadzhev

NIGGG-BAS
Acad. G. Bonchev str., bl. 3
1113, Sofia, Bulgaria
ggadjev@geophys.bas.bg

Iliya Georgiev

Metro State University of Denver
Campus Box 38B, P.O.Box173362
80217-3362, Denver, Colorado, USA
gueorgil@msudenver.edu

Ivan Georgiev

Institute of Information and
Communication Technologies &
Institute of Mathematics and Informatics
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 25A
1113 Sofia, Bulgaria
ivan.georgiev@parallel.bas.bg

Krassimir Georgiev

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 25A
1113 Sofia, Bulgaria
georgiev@parallel.bas.bg

Irina Georgieva

Institute of Mathematics and Informatics
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 8
1113 Sofia, Bulgaria
irina@math.bas.bg

Ivelina Georgieva

NIGGG-BAS
Acad. G. Bonchev str., bl. 3
1113, Sofia, Bulgaria
iivanova@geophys.bas.bg

Krassimira Georgieva

Metro State University of Denver
Campus Box 38B, P.O.Box173362
80217-3362, Denver, Colorado, USA

Silvia Grozdanova

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 25A
1113 Sofia, Bulgaria
silvia@parallel.bas.bg

Sylvi-Maria Gurova

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
and Faculty of Mathematics and
Informatics, Sofia University
Acad. G. Bonchev Str., bl. 25A
1113 Sofia, Bulgaria
smgurova@parallel.bas.bg

Stanislav Harizanov

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 25A
1113 Sofia, Bulgaria
sharizanov@parallel.bas.bg

Yanzhen Hou

Beijing Institute of Technology
5 South Zhongguancun Street
Haidian Zone
100081, Beijing, China
yanzhen@bit.edu.cn

Nevena Ilieva

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 25A
1113 Sofia, Bulgaria
nevena.ilieva@parallel.bas.bg

Aneta Karaivanova

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 25A
1113 Sofia, Bulgaria
anet@parallel.bas.bg

Kolyu Kolev

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 25A
1113 Sofia, Bulgaria
kkolev@iit.bas.bg

Mohamed Lachaab

University of Tunis
ISG
41 Rue de la liberte
2000, Bardo, Tunisia
mlachaab@albany.edu

Raytcho Lazarov

Texas A&M University
University Dr.
77843 Texas, USA
lazarov@math.tamu.edu

Elena Lilkova

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 25A
1113 Sofia, Bulgaria
elilkova@parallel.bas.bg

Angelos LIOLIOS

Democritus University of Thrace
L. Pyrgou 35 B
GR 67100 Xanthi, Greece
aliolios@civil.duth.gr

Asterios LIOLIOS

Democritus University of Thrace
L. Pirgou 35 B
GR 67100 Xanthi, Greece
liolios@civil.duth.gr

Konstantinos Liolios

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 25A
1113 Sofia, Bulgaria
kostisliolios@gmail.com

Ivan Lirkov

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 25A
1113 Sofia, Bulgaria
ivan@parallel.bas.bg

Svetozar Margenov

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 25A
1113 Sofia, Bulgaria
margenov@parallel.bas.bg

Maya Neytcheva

Uppsala University
Department of Information Technology
Box 337
75105, Uppsala, Sweden
maya.neytcheva@gmail.com

Dimitar Slavchev

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 25A
1113 Sofia, Bulgaria
dimitargslavchev@parallel.bas.bg

Jiri Stary

Institute of Geonics CAS
Studentska 1768
70800 Ostrava Poruba, Czech Republic
jiri.stary@ugn.cas.cz

Bálint Máté Takács

Eötvös Loránd University
Egyetem tér 1-3
1053, Budapest, Hungary
takacs.balint.mate@gmail.com

Petar Tomov

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 2
1113 Sofia, Bulgaria
p.tomov@iit.bas.bg

Ferdinand Evert Uilhoorn

Warsaw University of Technology
Gas Engineering Group
Nowowiejska 20
00-653, Warsaw, Poland
ferdinand.uilhoorn@pw.edu.pl

Yavor Vutov

Institute of Information and
Communication Technologies
Bulgarian Academy of Sciences
Acad. G. Bonchev Str., bl. 25A
1113 Sofia, Bulgaria
yavor.vutov@gmail.com

Zahari Zlatev

Aarhus University
Frederiksborgvej 399, P. O. Box 358
4000 Roskilde, Denmark
zz@envs.au.dk