# Parallel Computation of Sensitivity Analysis Data for the Danish Eulerian Model

Tzvetan Ostromsky[1], Ivan Dimov[1], Rayna Georgieva[1], and Zahari Zlatev[2]

[1] Institute of Information and Communication Technologies,
Bulgarian Academy of Sciences,
Acad. G. Bonchev str., bl. 25-A, 1113 Sofia, Bulgaria
{ceco,rayna}@parallel.bas.bg, ivdimov@bas.bg
http://www.bas.bg/iict/
[2] National Environmental Research Institute,
Department of Atmospheric Environment, Frederiksborgvej 399 P.O. Box 358,
DK-4000 Roskilde, Denmark
zz@dmu.dk
http://www.dmu.dk/AtmosphericEnvironment

**Abstract.** Sensitivity Analysis of the Danish Eulerian Model requires an extensive amount of output data from computationally expensive numerical experiments with a specially adapted for the purpose version of the model, called SA-DEM. It has been successfully implemented and run on the most powerful parallel supercomputer in Bulgaria - IBM Blue-Gene/P. A new enhanced version, capable of using efficiently the full capacity of the mashine, has recently been developed. It will be described in this paper together with some performance analysis and numerical results. The output results are used to construct some mesh-functions of ozone concentrations ratios to be used further in sensitivity analysis of the model by using Monte Carlo algorithms.

## 1 Introduction

The Unified Danish Eulerian Model (UNI-DEM) is a powerful air pollution model, used to calculate the concentrations of various dangerous pollutants and other species over a large geographical region (4800 × 4800 km), covering the whole of Europe, the Mediterranean and some parts of Asia and Africa. It takes into account the main physical, chemical and photochemical processes between the studied species, the emissions, the quickly changing meteorological conditions. This large and complex task is not suitable for direct numerical treatment. For the purpose of numerical solution it is split into submodels, which represent the main physical and chemical processes. The sequential splitting [5] is used in the production version of the model, although other splitting methods have also been considered and implemented in some experimental versions [1,4]. Spatial and time discretization makes each of the above submodels a huge computational task, challenging for the most powerful supercomputers available nowadays. That is why the parallelization has always been a key point in the computer implementation of DEM since its very early stages. A coarse-grain parallelization strategy

based on partitioning of the spatial domain appears to be the most efficient and well-balanced way on widest class of nowadays parallel machines (with not too many processors), although some restrictions apply. Other parallelizations are also possible and suitable to certain classes of supercomputers [7,8].

In the chemical submodel there is a number of parameters for control on the speed of the corresponding chemical reactions. By introducing some regular perturbations in these parameters we produce the necessary data to be used later in a new adaptive Monte Carlo approach to variance-based sensitivity analysis. In general, sensitivity analysis can help us to find out which simplifications can be done without significant loss of accuracy. It is also important to analyze the influence of variations of the chemical rate coefficients (and other parameters on a later stage), as there is always a certain level of uncertainty for their values. This knowledge can show us which parameters are most critical for a certain set of output results.

A special parallel version (SA-DEM) of the UNI-DEM has been created for this purpose [2,6] and efficiently implemented on the IBM BlueGene/P, the most powerful parallel machine in Bulgaria. Nevertheless, the sensitivity analysis task remains a huge computational problem, which requires enormous resources of storage and CPU time. Essential improvements of this version are made by introducing two new levels of parallelism (top-level(MPI) and bottom-level(OpenMP) respectively) in SA-DEM. They allow us to shorten many times the necessary computing time for obtaining the sensitivity analysis results and to use efficiently the IBM BlueGene/P machine up to its full capacity.

The general concept of sensitivity analysis and our utilization for the above problem is described briefly in Section 2. The mathematical background of Danish Eulerian Model and the scheme of its numerical solution are described in Section 3. Some details on parallelization of the improved SA-DEM version, performance and scalability results obtained on the IBM BlueGene/P are presented in the rest of this paper.

## 2   Sensitivity Analysis Concept — Sobol's Approach

Sensitivity analysis (SA) is the study of how much the uncertainty in the input data of a model (due to any reason: inaccurate measurements or calculation, approximation, data compression, etc.) is reflected in the accuracy of the output results [9]. Two kinds of sensitivity analysis are present in the existing literature, local and global. Local SA studies how much some small variations of inputs around a given value can change the value of the output. Global SA takes into account all the variation range of the input parameters, and apportions the output uncertainty to the uncertainty in the input data. Subject to our study in this paper is the global sensitivity analysis.

Several sensitivity analysis techniques have been developed and used throughout the years [9]. In general, these methods rely heavily on special assumptions connected to the behaviour of the model (such as linearity, monotonicity and additivity of the relationship between input and output parameters of the model).

Among the quantitative methods, variance-based methods are most often used. The main idea of these methods is to evaluate how the variance of an input or a group of inputs contributes to the variance of model output.

Assume that a model is represented by the following model function: $u = f(\mathbf{x})$, where the input parameters $\mathbf{x} = (x_1, x_2, \ldots, x_d) \in U^d \equiv [0,1]^d$ are independent (non-correlated) random variables with a known *joint probability distribution function*. In this way the output $u$ becomes also a random variable (as it is a function of the random vector $\mathbf{x}$) and let $\mathbf{E}$ be its mathematical expectation. Let $\mathbf{D}[\mathbf{E}(u|x_i)]$ be the variance of the conditional expectation of $u$ with respect to $x_i$ and $\mathbf{D}_u$ - the total variance according to $u$. This indicator is called *first-order sensitivity index* by Sobol [10] or sometimes *correlation ratio*.

*Total Sensitivity Index (TSI)* [10] of an input parameter $x_i$, $\;i \in \{1, \ldots, d\}$ is the sum of the complete set of mutual sensitivity indices of any order (main effect, two-way interactions (second order), three-way interactions (third order), etc.):

$$S_{x_i}^{tot} = S_i + \sum_{l_1 \neq i} S_{il_1} + \sum_{l_1, l_2 \neq i, l_1 < l_2} S_{il_1 l_2} + \ldots + S_{il_1 \ldots l_{d-1}}, \tag{1}$$

where $S_{il_1 \ldots l_{j-1}}$ – $j^{\text{th}}$ order sensitivity index for the parameter $\;x_i\;(1 \leq j \leq d)$, $j = 1:\;\;S_i$ – the "main effect" of $\;x_i$. According to the values of their total sensitivity indices, the input parameters are classified in the following way: very important $(0.8 < S_{x_i}^{tot})$, important $(0.5 < S_{x_i}^{tot} < 0.8)$, unimportant $(0.3 < S_{x_i}^{tot} < 0.5)$, irrelevant $(S_{x_i}^{tot} < 0.3)$. In most practical problems the high dimensional terms can be neglected, thus reducing significantly the number of summands in (1).

The Sobol's method is one of the most often used variance-based methods. It is based on a unique decomposition of the model function into orthogonal terms (summands) of increasing dimension and zero means. Its main advantage is computing in a uniform way not only the first order indices, but also the higher order indices (in quite a similar way as the computation of the main effects). The total sensitivity index can then be calculated with just one Monte Carlo integral per factor.

The Sobol's method for global SA, applied here, is based on the so-called $HDMR^1$ (2) of the model function $f$ (integrable) in the $d$-dimensional factor space:

$$f(\mathbf{x}) = f_0 + \sum_{s=1}^{d} \sum_{l_1 < \ldots < l_s} f_{l_1 \ldots l_s}(x_{l_1}, x_{l_2}, \ldots, x_{l_s}), \tag{2}$$

where $f_0$ is a constant. The representation (2) is not unique. Sobol has proven that under the conditions (3) for the right-hand-side functions

$$\int_0^1 f_{l_1 \ldots l_s}(x_{l_1}, x_{l_2}, \ldots, x_{l_s})\,\mathrm{d}x_{l_k} = 0, \quad 1 \leq k \leq s, \quad s = 1, \ldots, d \tag{3}$$

---

[1] High Dimensional Model Representation.

the decomposition (2) is unique and is called $ANOVA^2$-HDMR of the model function $f(\mathbf{x})$. Moreover, the functions of the right-hand side can be defined in a unique way by multidimensional integrals [11].

## 3   The Danish Eulerian Model

In this section we describe shortly the Danish Eulerian Model (DEM) [13] and its current production version UNI-DEM [12]. It is mathematically represented by the following system of partial differential equations, in which the unknown concentrations of a large number of chemical species (pollutants and other chemically active components) take part. The main physical and chemical processes (advection, diffusion, chemical reactions, emissions and deposition) are represented in that system.

$$
\begin{aligned}
\frac{\partial c_s}{\partial t} = & -\frac{\partial(uc_s)}{\partial x} - \frac{\partial(vc_s)}{\partial y} - \frac{\partial(wc_s)}{\partial z} + \\
& + \frac{\partial}{\partial x}\left(K_x\frac{\partial c_s}{\partial x}\right) + \frac{\partial}{\partial y}\left(K_y\frac{\partial c_s}{\partial y}\right) + \frac{\partial}{\partial z}\left(K_z\frac{\partial c_s}{\partial z}\right) + \\
& + E_s + Q_s(c_1, c_2, \ldots c_q) - (k_{1s} + k_{2s})c_s, \quad s = 1, 2, \ldots q\ .
\end{aligned}
\tag{4}
$$

where

- $c_s$ – the concentrations of the chemical species;
- $u$, $v$, $w$ – the wind components along the coordinate axes;
- $K_x$, $K_y$, $K_z$ – diffusion coefficients;
- $E_s$ – the emissions;
- $k_{1s}$, $k_{2s}$ – dry / wet deposition coefficients;
- $Q_s(c_1, c_2, \ldots c_q)$ – non-linear functions describing the chemical reactions between species under consideration.

The above rather complex system (4) is split (by using the most straightforward sequential splitting scheme) according to the major physical and chemical processes. Finaly, the following 3 submodels are formed:

$$
\frac{\partial c_s^{(1)}}{\partial t} = -\frac{\partial(uc_s^{(1)})}{\partial x} - \frac{\partial(vc_s^{(1)})}{\partial y} + \frac{\partial}{\partial x}\left(K_x\frac{\partial c_s^{(1)}}{\partial x}\right) + \frac{\partial}{\partial y}\left(K_y\frac{\partial c_s^{(1)}}{\partial y}\right) = A_1 c_s^{(1)}(t)
$$

**horizontal advection & diffusion**

$$
\frac{\partial c_s^{(2)}}{\partial t} = E_s + Q_s(c_1^{(2)}, c_2^{(2)}, \ldots c_q^{(2)}) - (k_{1s} + k_{2s})c_s^{(2)} = A_2 c_s^{(2)}(t)
$$

**chemistry, emissions & deposition**

$$
\frac{\partial c_s^{(3)}}{\partial t} = -\frac{\partial(wc_s^{(3)})}{\partial z} + \frac{\partial}{\partial z}\left(K_z\frac{\partial c_s^{(3)}}{\partial z}\right) = A_3 c_s^{(3)}(t)
$$

**vertical transport**

---

[2] ANalysis Of VAriances.

Spatial and time discretization of the above submodels on the EMEP[3] grid or its refinements (see Table 1) makes each of them a huge computational task. Thus the high performance and parallel computing become vital for the real-time numerical solution of the model.

The following methods are used in the numerical solution of the submodels:

- **Advection-diffusion part:** Finite elements, followed by predictor-corrector schemes with several different correctors.
- **Chemistry-deposition part:** An improved version of the QSSA (Quazi Steady-State Approximation)
- **Vertical transport:** Finite elements, followed by theta-methods.

## 4 UNI-DEM, the Improved Sensitivity Analysis Version SA-DEM and Their Parallel Implementation Features

The development and improvements of DEM throughout the years has lead to a variety of different versions with respect to the grid-size/resolution, vertical layering (2D or 3D model respectively) and the number of species in the chemical scheme. The most prospective of them have been united in the packege UNI-DEM. The available up-to-date versions, the selecting parameters and their optional values are shown in Table 1.

A coarse-grain parallelization strategy based on partitioning of the spatial domain in strips or blocks is currently used in UNI-DEM. For the purpose of this study, the strip-based distributed memory parallelization of the model via MPI is used [3,7,14]. It is based on partitioning of the horizontal grid, which implies certain restrictions on the number of MPI tasks and requires communication on each time step. Improving the data locality for more efficient cache utilization is achieved by using *chunks* to group properly the small tasks in the chemistry-deposition and vertical exchange stages. Additional pre-processing and post-processing stages are needed for scattering the input data and gathering the results, causing some overhead.

SA-DEM is a modification of UNI-DEM, specially adjusted to be used in the first stage of our sensitivity analysis concept (see [2]). There are additional input

**Table 1.** User-determined parameters for selecting an appropriate UNI-DEM version

| Parameter | Description | Optional values | | |
|-----------|-------------|-----------------|---|---|
| NX = NY | Grid size (Grid step) | $96 \times 96$ (50 km) | $288 \times 288$ (16.7 km) | $480 \times 480$ (10 km) |
| NZ | # layers (2D/3D) | 1 or 10 | | |
| NEQUAT | # chem. species | 35, 56 or 168 | | |

---

[3] European Monitoring and Evaluation Programme.

parameters in the main program, allowing the user to set some changes of the parameters subject to sensitivity analysis deeply in the code. These are constants in the original model and normaly there is no direct user access to their values. In our particular sensitivity analysis study regular perturbations have to be done on some chemical rate coefficients in the chemistry submodel. These coefficients must be modified on the course of the SA experiments, either separately or in groups in dependence with the dimension of the particular sensitivity analysis study. That is a typical SIMD[4] task, if considering the coarsest possible level of the strusture of our algorithm. By using it we introduce a new, higher level of parallelism in SA-DEM on the top of the grid-partitioning level, the basis for distributed-memory MPI parallelization in UNI-DEM.

Our target hardware can optionally offer a limited amount of shared memory parallelism. In order to exploit it efficiently, we introduced an additional (finer-grain) level of parallelism in our algorithm by using OpenMP standard directives.

All three levels of parallelism can be used efficiently in the calculations of the necessary data for sensitivity analysis on a powerful Blue Gene/P computing system. This is shown by experiments in the next section. Finally, for extracting the ozone mean monthly concentrations and computing the necessary mesh functions an additional program was developed. The last task is much simpler and not computationally intensive, so currently we left it beyond the scope of our highly parallel supercomputer implementation.

## 5    Numerical Experiments on the IBM Blue Gene/P

In this section we present some execution times and speed-ups in order to show the scalability of SA-DEM on the Bulgarian IBM Blue Gene/P , the main computing platform used in our sensitivity analysis study. The IBM Blue Gene/P is a state-of-the-art high-performance system with 8192 CPU in total and theoretical peak performance more than 23 TFLOPS. It consists of 2048 compute cards (nodes), each of them being a quard core PowerPC 450 (4 CPU, 850 MHz, 2 GB RAM). A single compute card is in fact a 4-CPU shared-memory computational unit with possible multithreading support via OpenMP. It can be used in 3 different modes: VN, DUAL and SMP. With respect to the MPI parallelism there are 4 MPI processes per node in VN mode, 2 - in DUAL mode, and one in VN mode. Thus, in the last two cases the machine offers limited, but natural from hardware viewpoint shared memory parallelism, exploited on the lowest (finer-grain) level in the new implementation of SA-DEM, as mentioned above. There is 8 MB L3 cache per node, 32 KB L1 cache per CPU (private).

The results of 20-sample one-year experiments with the SA-DEM (on the 2D medium resolution spatial grid (96 x 96 x 1)), executed on the Blue Gene/P are presented in Table 2 below.

The load managing policy of this huge parallel system is based on allocating whole number of planes per job (a multiple of 128 nodes). Therefore it does not encourage submission of long jobs that use considerably less nodes, as this

---

[4] Single Instruction Multiple Data, according to Flynn's taxonomy (1966).

**Table 2.** Time (T) in seconds and **speed-up (Sp)** of SA-DEM with MPI parallelism on the Bulgarian IBM Blue Gene/P (in VN mode)

| # | Advection | | Chemistry | | Comm. | I/O | TOTAL | | |
|---|---|---|---|---|---|---|---|---|---|
| CPU | time [s] | **(Sp)** | time [s] | **(Sp)** | time [s] | time [s] | time [s] | **(Sp)** | E [%] |
| 40 | 3410 | ( **40**) | 15925 | ( **40**) | 94 | 1116 | 20733 | ( **40**) | 100% |
| 80 | 1715 | ( **79**) | 7948 | ( **80**) | 99 | 1151 | 11000 | ( **75**) | 94% |
| 120 | 1154 | (**118**) | 5291 | (**120**) | 138 | 1051 | 7664 | (**108**) | 90% |
| 160 | 870 | (**157**) | 3983 | (**160**) | 137 | 1076 | 6204 | (**134**) | 84% |
| 240 | 586 | (**233**) | 2643 | (**241**) | 140 | 1107 | 4562 | (**182**) | 76% |
| 320 | 464 | (**294**) | 1974 | (**323**) | 153 | 1131 | 3810 | (**218**) | 68% |
| 480 | 344 | (**396**) | 1321 | (**482**) | 221 | 1651 | 3659 | (**227**) | 47% |
| 640 | 283 | (**482**) | 985 | (**647**) | 176 | 1973 | 3473 | (**239**) | 37% |
| 960 | 206 | (**662**) | 656 | (**971**) | 172 | 1972 | 3114 | (**266**) | 28% |

*Title within table:* Time and **speed-up** of SA-DEM on the IBM Blue Gene/P ($96 \times 96 \times 1$) grid,    35 species,    CHUNKSIZE=48

**Table 3.** Time (T) in seconds and **speed-up (Sp)** of SA-DEM with both MPI and OpenMP parallelism on the Bulgarian IBM Blue Gene/P

*Title within table:* Time and **speed-up** of SA-DEM (MPI+OpenMP) on the IBM Blue Gene/P ($96 \times 96 \times 1$) grid,    35 species,    CHUNKSIZE=48

| # | MPI p-s × | | Advection | | Chemistry | | TOTAL | | |
|---|---|---|---|---|---|---|---|---|---|
| CPU | OMP thr. | MODE | T [s] | **(Sp)** | T [s] | **(Sp)** | T [s] | **(Sp)** | E [%] |
| 40 | 40 × 1 | VN | 3410 | ( **40**) | 15925 | ( **40**) | 20733 | ( **40**) | 100% |
| 80 | 40 × 2 | DUAL | 1778 | ( **77**) | 7972 | ( **80**) | 11295 | ( **73**) | 92% |
| 160 | 80 × 2 | DUAL | 889 | (**153**) | 3960 | (**161**) | 6153 | (**135**) | 84% |
| 240 | 120 × 2 | DUAL | 647 | (**211**) | 2655 | (**240**) | 4712 | (**176**) | 73% |
| 320 | 160 × 2 | DUAL | 502 | (**271**) | 1978 | (**322**) | 4006 | (**207**) | 65% |
| 480 | 240 × 2 | DUAL | 358 | (**381**) | 1329 | (**479**) | 3418 | (**243**) | 51% |
| 640 | 160 × 4 | SMP | 223 | (**612**) | 997 | (**639**) | 2768 | (**300**) | 47% |
| 960 | 480 × 2 | DUAL | 218 | (**626**) | 659 | (**967**) | 2684 | (**309**) | 32% |
| 960 | 240 × 4 | SMP | 158 | (**863**) | 667 | (**955**) | 2292 | (**362**) | 38% |
| 1280 | 320 × 4 | SMP | 122 | (**1118**) | 499 | (**1277**) | 2109 | (**393**) | 31% |
| 1920 | 480 × 4 | SMP | 99 | (**1378**) | 338 | (**1885**) | 2568 | (**323**) | 17% |
| 2560 | 640 × 4 | SMP | 83 | (**1643**) | 332 | (**1919**) | 2182 | (**380**) | 15% |
| 3840 | 960 × 4 | SMP | 58 | (**2352**) | 168 | (**3792**) | 1653 | (**502**) | 13% |

would be a waste of resourse. In our experiments the runs start from 40 CPU. In order to obtain comparable figures and correct scalability results we calculate the speed-up in the next tables by using the following formula (assuming that the speed-up on 40 processors is 40):

$$Sp(n) = 40 \, \frac{T(n)}{T(40)} \qquad (5)$$

where $n$ is the number of processors (given in the first column). The time and the **speed-up (Sp)** of the main computational stages and in total are given in separate columns. The last column contains also the total efficiency $E$ (in percent), where $E = 100 \, Sp(n)/n \, \%$.

The total time includes also the MPI communication time as well as the time for some I/O procedures, which are not parallelizable. Moreover, the larger the number of MPI tasks, the more I/O device conflicts arise, which results in a significant drop-down in the total efficiency. I/O device access appear to be the performance bottleneck in this case, partially avoided by using the lowest level OpenMP parallelisation (see Table 3). On the other hand, the computational stages scale pretty well, even the speed-up of the chemistry stage tends to be slightly superlinear (due to the cache memory effects).

## 6 Conclusions and Plans for Future Work

We consider a 3-stage variance-based sensitivity analysis method. For the purpose of sensitivity analysis of the Danish Eulerian Model with respect to variation of certain chemical rate coefficients, a special version of the model has been developed and implemented efficiently on the IBM Blue Gene/P (called SA-DEM). Experiments, showing its scalability and efficiensy on a huge parallel system (IBM Blue Gene/P) are presented in this paper.

The first stage of our 3-stage sensitivity analysis method is completed by extracting from the output results some mean monthly concentrations of the ozone and producing the necessary mesh functions. The second stage of this sensitivity analysis research includes approximation of the mesh functions by polynomials of 3-rd / 4-th degree or by cubic B-spline functions. A Monte Carlo integration method is furtherly applied to these functions on the third stage. The results of the last two stages are presented in another paper.

Our near future plans include:

- Optimization of the I/O operations in order to overcome the bottleneck, causing a significant efficiency dropdown;
- Extending the abilities of SA-DEM (including experiments with more chemical species and on finer resolution grids (storage-permitting);
- Extending the scope of the sensitivity analysis study with respect to the emission levels and the boundary conditions.

# References

1. Dimov, I., Faragó, I., Havasi, Á., Zlatev, Z.: Operator splitting and commutativity analysis in the Danish Eulerian Model. Math. Comp. Sim. 67, 217–233 (2004)
2. Dimov, I., Georgieva, R., Ivanovska, S., Ostromsky, T., Zlatev, Z.: Studying the Sensitivity of the Pollutants Concentrations Caused by Variations of Chemical Rates. Journal of Computational and Applied Mathematics 235(2), 391–402 (2010)
3. Dimov, I., Georgiev, K., Ostromsky, T., Zlatev, Z.: Computational challenges in the numerical treatment of large air pollution models. Ecological Modelling 179, 187–203 (2004)
4. Dimov, I., Ostromsky, T., Zlatev, Z.: Challenges in using splitting techniques for large-scale environmental modeling. In: Faragó, I., Georgiev, K., Havasi, Á. (eds.) Advances in Air Pollution Modeling for Environmental Security. NATO Science Series, vol. 54, pp. 115–132. Springer, Heidelberg (2005)
5. Marchuk, G.I.: Mathematical modeling for the problem of the environment. Studies in Mathematics and Applications, vol. 16. North-Holland, Amsterdam (1985)
6. Ostromsky, T., Dimov, I., Georgieva, R., Zlatev, Z.: Sensitivity Analysis of a Large-Scale Air Pollution Model: Numerical Aspects and a Highly Parallel Implementation. In: Lirkov, I., Margenov, S., Waśniewski, J. (eds.) LSSC 2009. LNCS, vol. 5910, pp. 197–205. Springer, Heidelberg (2010)
7. Ostromsky, T., Zlatev, Z.: Parallel Implementation of a Large-Scale 3-D Air Pollution Model. In: Margenov, S., Waśniewski, J., Yalamov, P. (eds.) LSSC 2001. LNCS, vol. 2179, pp. 309–316. Springer, Heidelberg (2001)
8. Ostromsky, T., Zlatev, Z.: Flexible Two-Level Parallel Implementations of a Large Air Pollution Model. In: Dimov, I.T., Lirkov, I., Margenov, S., Zlatev, Z. (eds.) NMA 2002. LNCS, vol. 2542, pp. 545–554. Springer, Heidelberg (2003)
9. Saltelli, A., Tarantola, S., Campolongo, F., Ratto, M.: Sensitivity Analysis in Practice: A Guide to Assessing Scientific Models. Halsted Press, New York (2004)
10. Sobol, I.M.: Sensitivity estimates for nonlinear mathematical models. Mathematical Modeling and Computational Experiment 1, 407–414 (1993)
11. Sobol, I.M.: Global Sensitivity Indices for Nonlinear Mathematical Models and Their Monte Carlo Estimates. Mathematics and Computers in Simulation 55(1-3), 271–280 (2001)
12. WEB-site of the Danish Eulerian Model,
    `http://www.dmu.dk/AtmosphericEnvironment/DEM`
13. Zlatev, Z.: Computer treatment of large air pollution models. Kluwer (1995)
14. Zlatev, Z., Dimov, I., Georgiev, K.: Three-dimensional version of the Danish Eulerian Model. Zeitschrift für Angewandte Mathematik und Mechanik 76(S4), 473–476 (1996)