

## THE MUST MODEL EVALUATION EXERCISE: STATISTICAL ANALYSIS OF MODELLING RESULTS

J. Franke<sup>1</sup>, J. Bartzis<sup>2</sup>, F. Barmpas<sup>3</sup>, R. Berkowicz<sup>4</sup>, K. Brzozowski<sup>5</sup>, R. Buccolieri<sup>6</sup>, B. Carissimo<sup>7</sup>,  
A. Costa<sup>8</sup>, S. Di Sabatino<sup>6</sup>, G. Efthimiou<sup>2</sup>, I. Goricsan<sup>9</sup>, A. Hellsten<sup>10</sup>, M. Ketzel<sup>1</sup>, B. Leidl<sup>11</sup>,  
R. Nuterman<sup>12</sup>, H. Olesen<sup>4</sup>, E. Polreich<sup>13</sup>, J. Santiago<sup>14</sup>, R. Tavares<sup>8</sup>

<sup>1</sup>Department of Fluid- and Thermodynamics, University of Siegen, Germany

<sup>2</sup>Univ. of West Macedonia, Greece; <sup>3</sup>Aristotle Univ., Greece; <sup>4</sup>NERI, Univ. of Aarhus, Denmark;

<sup>5</sup>Univ. of Bielsko-Biala, Poland; <sup>6</sup>Univ. of Salento, Italy; <sup>7</sup>EDF, France; <sup>8</sup>Univ. of Aveiro, Portugal;

<sup>9</sup>Budapest Univ. of Technology and Economics, Hungary; <sup>10</sup>Helsinki Univ. of Technology, Finland;

<sup>11</sup>Univ. of Hamburg, Germany; <sup>12</sup>Tomsk State Univ., Russia; <sup>13</sup>ZAMG, Austria; <sup>14</sup>CIEMAT, Spain

**Abstract:** The first validation exercise of the COST action 732 lead to a substantial number of simulation results for comparison with the MUST wind tunnel experiments. Validation metrics for selected simulation results of the flow field and the concentrations are presented and compared to the state of the art. In addition mean metrics and corresponding scatter limits are computed from the individual results.

**Key words:** COST 732, MUST, CFD, atmospheric dispersion, model evaluation, model validation, metrics, N-Version testing.

### 1. INTRODUCTION

The increasing use of computational simulation for dispersion predictions at micro-scale aggravates the need for a clear and commonly accepted procedure for the quality assurance of the codes and the simulation results. To that end the European COST action 732 *Quality Assurance and Improvement of Micro-Scale Meteorological Models* has proposed a protocol for the evaluation of codes (Britter and Schatzmann, 2007). This protocol has been applied by the action's participants to the mocked urban setting test (MUST) case. For this case both field (Yee and Bilitoft, 2004) and wind tunnel measurements (Leidl et al., 2007) are available. In the first step the validation of the codes was done against wind tunnel measurements of velocities and concentrations for several wind directions.

The validation part of the protocol recommends exploratory data analysis and statistical performance measures to assess the predictive capability of the codes. First results of the exploratory data analysis are presented by Olesen et al. (2008) in a companion paper. Here selected results for the statistical performance measures are presented. These measures are defined as validation metrics, allowing a quick assessment of the simulation quality by delivering one value for a large number of measurement points.

Furthermore the metrics are analysed statistically. This is possible because of the large number of simulation results that were produced by the 13 different research groups participating in the exercise. To that end the scatter of the different results is viewed as reproducibility of the computational process. This is known as N-Version testing in computation and equals the N-th order replication in experiments (Coleman and Steele, 1999). The outcome is a collective metric of all simulations, together with scatter limits, which can also be used for the certification of CFD codes (Stern et al., 2006).

### 2. EXPERIMENTAL SETUP AND DEFINITION OF STATISTICS

#### Measurement positions

Extensive wind tunnel measurements of the flow field and the concentrations were conducted for several wind directions (Bezplacova, 2007; Harms et al., 2005). In this work only velocity measurements at 18 towers and dispersion measurements for the -45 degree approach flow case are used for the calculation of metrics, as these data are also used in the exploratory data analysis of Olesen et al. (2008). In Figure 1 the tower positions and the concentration measurement positions are shown. Altogether there are 498 Velocity measurement positions and 256 concentration measurement positions.

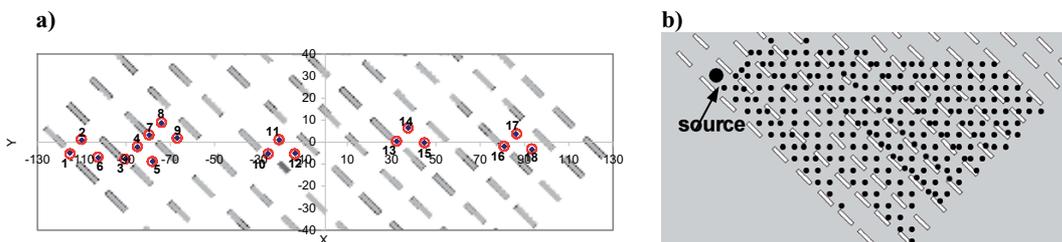


Figure 1. Measurement positions for the -45 degree case. a) towers for velocities, b) concentrations. Wind blows from the left.

### Validation metrics

For the velocity data and the concentrations the hit rate  $q$  is proposed as one metric (VDI, 2005).

$$q = \frac{1}{I} \sum_{n=1}^I i_n \quad \text{with} \quad i_n = \begin{cases} 1 & \text{if } |(O_n - P_n)/O_n| \leq \Delta_r \text{ or } |O_n - P_n| \leq \Delta_a \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

Here  $N$  = number of measurement positions;  $O$  = observed value;  $P$  = predicted value;  $\Delta_r = 0.25$ , the allowed relative difference;  $\Delta_a$  = allowed absolute difference. For the allowed absolute difference the measurement uncertainties for the corresponding variables are used, see Table 1.

Table 1. Allowed absolute differences  $\Delta_a$  used in the calculation of the hit rate, Equation (1).

| $U/U_{ref}$ | $W/U_{ref}$ | $k/U_{ref}^2$ | $C^*$ |
|-------------|-------------|---------------|-------|
| 0.008       | 0.007       | 0.005         | 0.003 |

Here  $C^*$  is the scaled concentration

$$C^* = C \cdot U_{ref} \cdot H^2 / Q_{source} \quad (2)$$

where  $U_{ref}$  = reference velocity in  $x$ -direction at  $(x, y, z) = (-144, -2.25, 7.29)$  m and  $Q_{source}$  = volumetric flow rate of the source, see Figure 1b.

For the concentrations in addition the metrics also used in the BOOT software (Chang and Hanna, 2004) are calculated from the observed (index o) and predicted (index p) concentrations  $C^*$ . These are the factor of two  $FAC2$ , the fractional bias  $FB$ , the normalised mean square error  $NMSE$ , the geometric mean  $MG$  and the geometric variance  $VG$ , which are defined as

$$\begin{aligned} FAC2 &= \text{fraction of data with } 0.5 \leq C_p^* / C_o^* \leq 2, \\ FB &= 2 \left( \langle C_o^* \rangle - \langle C_p^* \rangle \right) / \left( \langle C_o^* \rangle + \langle C_p^* \rangle \right), \quad NMSE = \left\langle \left( C_o^* - C_p^* \right)^2 \right\rangle / \left( \langle C_o^* \rangle \cdot \langle C_p^* \rangle \right), \\ MG &= \exp \left( \langle \ln C_o^* \rangle - \langle \ln C_p^* \rangle \right), \quad VG = \exp \left[ \left\langle \left( \ln C_o^* - \ln C_p^* \right)^2 \right\rangle \right], \end{aligned} \quad (3)$$

where angular brackets denote an average over all measurement points. For  $MG$  and  $VG$  the experimental uncertainty from Table 1 is used as threshold, i.e. for the observed and predicted values  $\max(0.003, C^*)$  is used. In the calculation of  $FAC2$  it is checked for low concentrations if both, observed and predicted are below the threshold. If so, then the position is within a factor of two. All metric calculations are performed within the Excel workbooks described by Olesen et al. (2008).

### Statistics of metrics

Assuming that all simulation results belong to a Gaussian or normal parent population the mean  $Y$  and the standard deviation  $S$  of the sample population are defined as

$$Y = 1/M \sum_{j=1}^M X_j, \quad S = \left[ 1/(M-1) \sum_{j=1}^M (X_j - Y)^2 \right]^{1/2} \quad (4)$$

where  $X_j$  = one of the metrics presented above for simulation  $j$  and  $M = 20$  is the number of simulation results. As this definition of the mean is known to be easily contaminated by outliers (Müller, 2000), additionally the median  $Z$  of the sample population is computed and the standard deviation  $T$  based on the sample median absolute deviation (MAD), see e.g. Hemsch (2000).

$$T = 1/0.6745 \sqrt{M/(M-1)} \text{median}(|X_i - Z|) \quad (5)$$

From the standard deviations scatter limits are computed as  $P_S = 2.093S$  and  $P_T = 2.093T$ , corresponding to 95% confidence intervals around the mean,  $Y \pm P_S$  and  $Z \pm P_T$  (Coleman and Steele, 1999).

## 3. RESULTS

Up to 30 simulation results from 13 different research groups are at present available. The statistical results are shown for 20 of these simulation results, which are all obtained with prognostic models, comprising general purpose commercial CFD codes, commercial micro-scale CFD codes and in house CFD or micro-scale meteorological models. As a general overview of the predictive capabilities shall be given the results are presented anonymously and no ranking of the codes is performed.

### Metrics for the flow field at the towers of the -45 degree case

In Figure 2 the hit rates for the two measured velocity components at the 18 towers are shown. Also included is the lower limit of  $q = 0.66$  as broken line, which is used by the VDI (2005) for its test cases to define successful validation.

While the velocity component in flow direction,  $U/U_{ref}$  is predicted very well by nearly all simulations with hit rates larger than 0.66, the vertical velocity component,  $W/U_{ref}$  is not well predicted by all simulations. Olesen et al. (2008)

showed that at the towers  $W/U_{ref}$  is mostly underpredicted. This shortcoming is however well reflected in the low hit rates. The quickly available metrics are therefore useful for identifying possible problems that can then be further analysed by an detailed inspection of the simulation results at the measurement positions. This does however not mean that a detailed analysis of the data is obsolete in the case of good metrics. E.g. the high rates of  $U/U_{ref}$  are mostly due to the expected good agreement between simulation and experiment well above the containers, while below container height the agreement deteriorates.

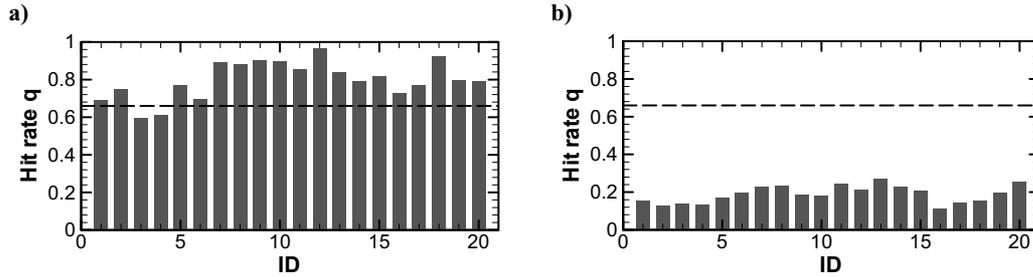


Figure 2. Hit rate for mean velocity components at the towers for the  $-45$  degree case. a)  $U/U_{ref}$  b)  $W/U_{ref}$ .

The hit rate for the turbulent kinetic energy  $k/U_{ref}^2$  at all towers is shown in Figure 3a. Here large differences between the different simulations are visible. As  $k/U_{ref}^2$  is a positive quantity also other metrics which are normally only used for concentrations are meaningful and can be analysed. E.g. the large positive  $FB$ s in Figure 3b clearly indicate that those simulations with low hit rates underpredict the turbulent kinetic energy.

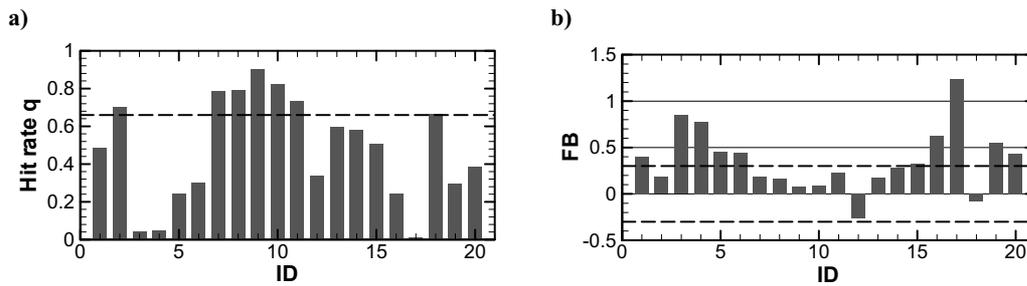


Figure 3. Metrics for the turbulent kinetic energy  $k/U_{ref}^2$  at the towers for the  $-45$  degree case. a) Hit rate  $q$ , b)  $FB$ .

### Metrics for the concentrations of the $-45$ degree case

The metrics of dispersion results are shown in Table 2. Numbers in italics are results which are not within the range defined by Chang, J.C. and S.R. Hanna (2004) as state of the art,  $FAC2 > 0.5$ ,  $|FB| < 0.3$ ,  $NMSE < 4$ ,  $0.7 < MG < 1.3$  and  $VG < 1.6$ .

Table 2. Dispersion metrics for the  $-45$  degree case. Italics indicate results being not state of the art.

| Simulation ID | Hit rate $q$ | FAC2 | FB           | NMSE         | MG          | VG          |
|---------------|--------------|------|--------------|--------------|-------------|-------------|
| 1             | <i>0.46</i>  | 0.71 | <i>-0.54</i> | <i>7.52</i>  | 0.82        | <i>1.65</i> |
| 2             | <i>0.65</i>  | 0.87 | <i>-0.36</i> | <i>7.07</i>  | 0.93        | 1.23        |
| 3             | <i>0.39</i>  | 0.65 | <i>-0.91</i> | <i>20.39</i> | <i>0.69</i> | <i>1.79</i> |
| 4             | <i>0.37</i>  | 0.67 | <i>-0.82</i> | <i>13.08</i> | 0.72        | <i>1.66</i> |
| 5             | <i>0.43</i>  | 0.64 | <i>-0.50</i> | <i>8.80</i>  | 0.94        | <i>1.67</i> |
| 6             | 0.80         | 0.88 | <i>-0.33</i> | <i>10.34</i> | 1.00        | 1.26        |
| 7             | 0.74         | 0.86 | -0.30        | <i>4.52</i>  | 0.86        | 1.23        |
| 8             | <i>0.57</i>  | 0.89 | <i>-0.38</i> | 2.72         | 0.84        | 1.20        |
| 9             | 0.68         | 0.97 | <i>-0.27</i> | 1.65         | 0.86        | 1.10        |
| 10            | <i>0.42</i>  | 0.58 | <i>0.40</i>  | <i>4.19</i>  | <i>1.36</i> | <i>1.70</i> |
| 11            | <i>0.54</i>  | 0.83 | <i>0.41</i>  | 3.54         | <i>1.42</i> | 1.36        |
| 12            | <i>0.63</i>  | 0.79 | 0.27         | 1.76         | <i>1.34</i> | 1.43        |
| 13            | <i>0.48</i>  | 0.62 | <i>-0.57</i> | <i>8.49</i>  | 0.89        | <i>1.72</i> |
| 14            | <i>0.38</i>  | 0.52 | <i>-0.62</i> | <i>12.51</i> | 0.86        | <i>2.27</i> |
| 15            | 0.79         | 0.86 | <i>-0.41</i> | <i>9.83</i>  | 0.81        | 1.41        |
| 16            | <i>0.58</i>  | 0.79 | 0.15         | 2.89         | 1.15        | 1.37        |
| 17            | <i>0.64</i>  | 0.82 | <i>-0.36</i> | <i>4.56</i>  | 0.89        | 1.39        |
| 18            | <i>0.64</i>  | 0.95 | <i>-0.40</i> | 2.21         | 0.84        | 1.13        |
| 19            | <i>0.54</i>  | 0.73 | <i>-0.67</i> | <i>8.65</i>  | <i>0.68</i> | <i>2.05</i> |
| 20            | <i>0.52</i>  | 0.79 | <i>-0.43</i> | <i>4.94</i>  | 0.87        | 1.35        |

While there are only four simulations that have larger hit rates for the concentrations than 0.66, the limit defined by the VDI (2005), all simulations have  $FAC2 > 0.5$ . Looking at all simulation results there is no obvious correlation between the flow field results presented in the previous section and the concentration metrics. While the low hit rates of simulations 3 and 4 correspond to low hit rates for  $U/U_{ref}$  (Fig. 2a) and  $k/U_{ref}^2$  (Fig. 3a), simulation 17 has a rather high hit rate for the concentrations despite the extremely low hit rate for  $k/U_{ref}^2$ . The most probable reason for this behaviour is that errors in the computation of the flow field and errors in the prediction of the concentrations cancel each other, leading to rather good results for the concentration metrics. Taking all metrics into account is therefore necessary to identify situations where good results are predicted due to the wrong reasons.

From Table 2 it can be also seen that most of the simulations overpredict the concentrations as a negative  $FB$  corresponds to an overprediction. In general the concentrations close to the source (see Fig. 1b) are overpredicted. The large differences at these few positions lead to  $FB$ s with too high magnitudes. Therefore only four simulations have a  $FB$  within the range defined as state of the art. From these four results one has  $NMSE$  larger than 4, indicating that this simulation substantially overpredicts as well as underpredicts the observations. While the over- and underpredictions cancel in  $FB$ , they lead to a too high  $NMSE$ . As can be seen very clearly from Figure 4a) only three simulations therefore have both,  $FB$  and  $NMSE$  in the range defining the state of the art.

The geometric metrics  $MG$  and  $VG$  are not dominated by high concentration values. Therefore they do not suffer from the overpredictions close to the source. Only five simulations have a geometric mean  $MG$  outside the range defining state of the art. However, if also the results for the geometric variance are taken into account, there are only 10 simulations within the admissible range spanned by both metrics, see Figure 4b).

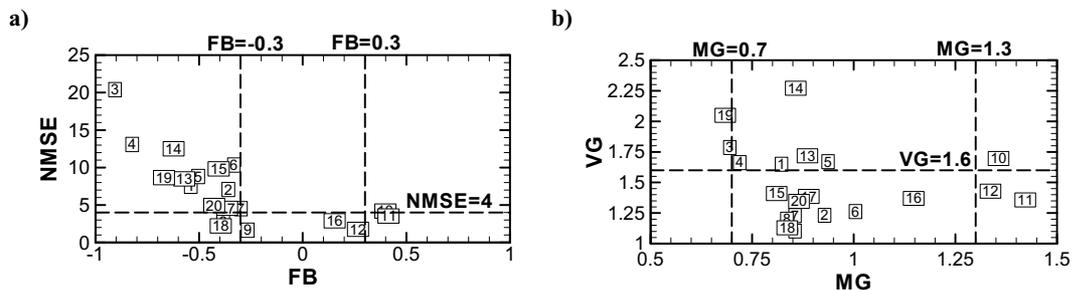


Figure 4. Metrics for the concentrations  $C^*$  of the  $-45$  degree case. a)  $NMSE$  vs.  $FB$ , b)  $VG$  vs.  $MG$ .

#### Statistics of dispersion metrics for the $-45$ degree case

In Figure 5 the running records of the hit rate  $q$  and the  $FB$  are shown together with the corresponding mean, median and associated scatter limits defining the 95% confidence interval. On the basis of statistics differences of the individual values within the confidence limits are considered as noise. Only those simulation results which lie outside the confidence intervals are then significant for further analysis.

For the hit rate the mean and the median are equal up to two digits,  $Y = Z = 0.56$ . The 95% confidence interval around the mean value is  $P_S = 0.28 = 0.5Y$ . For the median it is  $P_T = 0.31 = 0.55Z$ . Both intervals are large with regards to the mean value. All simulation results are within these intervals, there are no outliers. For  $FB$  the situation is different. Based on the median,  $Z = -0.39$ , and its confidence interval,  $P_T = 0.44 = 1.13|Z|$ , there are 5 outliers. The mean value,  $Y = -0.33$ , and its even larger confidence interval with  $P_S = 0.76 = 2.30|Y|$  do not have any outliers. This demonstrates the better suitability of the median to identify outliers which in the case of the mean value are only reflected by a larger confidence interval.

Except for  $FAC2$ ,  $MG$  and  $VG$  the collective metrics mean and median are all outside the limits defined as state of the art. This displays a systematic shortcoming of the models and indicates the need for at least an improved model set up, if not for model improvement.

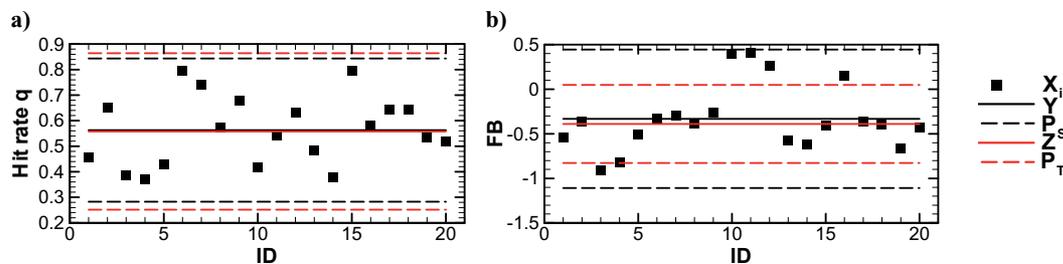


Figure 5. Running records of the metrics for the concentrations  $C^*$  of the  $-45$  degree case. a) Hit rate  $q$ , b)  $FB$ .

#### 4. CONCLUSIONS

Within the first validation exercise of the COST action 732 for the MUST wind tunnel experiment a large number of simulation results was produced by different research groups, using different models. The validation metrics for 20 of these results have been presented for one wind direction, comprising metrics for the flow field at selected locations and metrics for the concentrations at all concentration measurement positions. While all concentration results have a  $FAC2 > 0.5$ , other metrics are often outside the range defined as state of the art. In fact there is only one simulation result with all concentration metrics being state of the art.

This shortcoming of the individual results is also reflected in the collective results for the metrics, computed by statistical methods with the assumption that the 20 individual results are samples from a Gaussian parent population. From the collective metrics only  $FAC2$ ,  $MG$  and  $VG$  are within the range being state of the art. Another result from this statistical analysis is that the median as more robust estimate for the sample mean value, together with its standard deviation based on the median of the absolute deviations, is better suited for identifying outliers, i.e. individual simulations that differ substantially from the core region of the statistical distribution.

These results for the validation metrics indicate at least a need for an improved model set up for the MUST test case. After further detailed analysis of the individual simulation results recommendations for the conduction of simulations for the MUST case will be included in the official documents of the action, which can be found through URL 1.

#### REFERENCES

- Bezpalcova, K., 2007: Physical Modelling of Flow and Dispersion in an Urban Canopy, PhD thesis, Faculty of Mathematics and Physics, Charles University, Prague, 193 pp.
- Britter, R. and M. Schatzmann, 2007: Model evaluation guidance and protocol document, COST office, Belgium, 28 pp.
- Chang, J. C. and S.R. Hanna, 2004: Air quality model performance evaluation. *Meteo. Atmos. Phys.*, **87**, 167-196.
- Coleman, H.W. and W.G. Steele, 1999: Experimentation and Uncertainty Analysis for Engineers, 2<sup>nd</sup> Edition, John Wiley & Sons, USA, 275 pp.
- Harms, F., B. Leidl and M. Schatzmann, 2005: Comparison of tracer dispersion through a model of an idealized urban area from field (MUST) and wind tunnel measurements. *Proceedings International Workshop on Physical Modelling of Flow and Dispersion Phenomena*, London, Ontario, August 24-26.
- Leidl, B., K. Bezpalcova and F. Harms, 2007: Wind Tunnel Modelling Of The MUST Experiment. *11th Int. Conf. on Harmonisation within Atmospheric Dispersion Modelling for Regulatory Purposes*, Cambridge, July 2-5.
- Hensch, M., 2002: Statistical analysis of CFD solutions from the drag prediction workshop. *AIAA paper*, 2002-0842.
- Müller, J.W., 2000: Possible advantages of a robust evaluation of comparisons. *J. Res. Natl. Inst. Stand. Technol.*, **105**, 551-555.
- Olesen, H.R., A. Baklanov, J. Bartzis, F. Bampas, R. Berkowicz, K. Brzozowski, R. Buccolieri, B. Carissimo, A. Costa, S. Di Sabatino, G. Efthimiou, J. Franke, I. Goricsan, A. Hellsten, M. Ketzler, B. Leidl, R. Nuterman, E. Polreich, J. Santiago and R. Tavares, 2008: The MUST model evaluation exercise: Patterns in model performance. *12th Int. Conf. on Harmonisation within Atmospheric Dispersion Modelling for Regulatory Purposes*, Cavtat, October 6-9.
- Stern, F., R. Wilson and J. Shao, 2006: Quantitative V&V of CFD simulations and certifications of CFD code. *Int. J. Numer. Meth. Fluids*, **50**, 1335-1355.
- VDI, 2005: Environmental meteorology – Prognostic microscale windfield models – Evaluation for flow around buildings and obstacles. VDI guideline 3783, Part 9, Beuth, Berlin.
- Yee, E. and C.A. Biltoft, 2004: Concentration Fluctuation Measurements in a Plume Dispersing Through a Regular Array of Obstacles. *Boundary-Layer Meteorology*, **111**, 363-415.
- URL 1: Official web site of COST 732: <http://www.mi.uni-hamburg.de/index.php?id=464>