

Enhancing CFD-LES air pollution prediction accuracy using data assimilation

Elsa Aristodemou^{a,c}, Rossella Arcucci^{b,*}, Laetitia Mottet^{c,d}, Alan Robins^e,
Christopher Pain^{c,b}, Yi-Ke Guo^b

^a*School of Engineering, London South Bank University, UK*

^b*Data Science Institute, Department of Computing, Imperial College London, UK*

^c*Department of Earth Science & Engineering, Imperial College London, UK*

^d*Department of Architecture, University of Cambridge, UK*

^e*Department of Mechanical Engineering Sciences, University of Surrey, UK*

Abstract

It is recognised worldwide that air pollution is the cause of premature deaths daily, thus necessitating the development of more reliable and accurate numerical tools. The present study implements a three dimensional Variational (3DVar) data assimilation (DA) approach to reduce the discrepancy between predicted pollution concentrations based on Computational Fluid Dynamics (CFD) with the ones measured in a wind tunnel experiment. The methodology is implemented on a wind tunnel test case which represents a localised neighbourhood environment. The improved accuracy of the CFD simulation using DA is discussed in terms of absolute error, mean squared error and scatter plots for the pollution concentration. It is shown that the difference between CFD results and wind tunnel data, computed by the mean squared error, can be reduced by up to three order of magnitudes when using DA. This reduction in error is preserved in the CFD results and its benefit can be seen through several time steps after re-running the CFD simulation. Subsequently an optimal sensors positioning is proposed. There is a trade-off between the accuracy and the number of sensors. It was found that the accuracy was improved when placing/considering the sensors which were near the

*r.arcucci@imperial.ac.uk

pollution source or in regions where pollution concentrations were high. This demonstrated that only 14% of the wind tunnel data was needed, reducing the mean squared error by one order of magnitude.

Keywords: CFD, Data Assimilation, Wind Tunnel, Fluidity, Urban Environment, Pollutant concentration, Sensor positioning

1. Introduction

Climate change and air pollution form one of the grand challenges currently faced by humanity worldwide with still many questions remaining unanswered at the micro-scale/city-scale. The World Health Organisation (WHO) has found that outdoor air pollution in cities has been the primary cause of 4.2 million premature deaths annually worldwide [1]. WHO has subsequently established guidelines on the two main pollutants: the PM_{2.5} and NO₂ [1]. By 2030, reducing deaths and illnesses due to air pollution is also one of the aims of the United Nation sustainable development programme, with one of the goals being “good health and well being” [2]. In Europe, the European Union Commission has also already established procedures for monitoring and advising on air quality, focusing on the five most important air pollutants: ozone (O₃), nitrogen dioxide (NO₂), sulphur dioxide (SO₂), PM_{2.5} and PM₁₀ particles [3]. Currently, many areas in the United Kingdom (UK) and London in particular, fail to meet the WHO guidelines on two pollutants, the PM_{2.5} and NO₂. Due to these failings, the UK Government has developed the Clean Air Strategy 2019 [4], which sets out the UK plans for dealing with all sources of air pollution, and ensuring the health of the na-

19 tion through better air quality. In addition, for London, UK, the Mayor's
20 office has also developed the London Environment Strategy specifically for
21 air pollution problems in the capital [5].

22 It is therefore clear that serious steps are taken at both international, national
23 and city levels to reduce air pollution levels. Scientific and technological ad-
24 vances are therefore encouraged in the global effort to combat air pollution,
25 with innovative tools being developed to assist in this effort. Computational
26 methods/tools are at the forefront of these efforts, with many researchers
27 worldwide looking at how to most accurately capture the dispersion of pollu-
28 tants at the micro-scale level, within the urban environment [6, 7, 8]. Many
29 studies have been carried out over the years, employing both simplified Gaus-
30 sian plume models [8], to the more sophisticated ones using complex com-
31 putational fluid dynamics (CFD) with turbulence models ranging from the
32 Reynolds Averaged Navier-Stokes (RANS) approach to the more elaborate
33 Large Eddy Simulation (LES) methods [9, 10]. To validate the CFD sim-
34 ulations, comparison of various variables (velocity, pollutant concentration,
35 wind pressure coefficients, Reynolds stresses...) at micro-scale are usually
36 confronted to wind tunnel experiments data [11, 12, 13] and in less extend
37 to full scale experiments [14, 15]. For simple test case, i.e. simple geometric
38 configuration, CFD models are reliable and reproduce with a good agreement
39 data obtained from experiments [14]. However, the success of the comparison
40 becomes mitigated and significant discrepancies between CFD and experi-
41 ments are locally observed when more complex urban environment set-up

42 are considered [11, 12, 14]. In the context of predicting accurately the level
43 of pollution at pedestrian level, i.e. at micro-scale, more advance numerical
44 models need to be used in order to improve the reliability of their predictions.

45 The use of Data Assimilation (DA) technologies is a good candidate to an-
46 swer this need. DA is an uncertainty quantification technique used to in-
47 corporate observational data into a prediction model in order to improve
48 numerical forecasted results [16]. During the last 20 years, data assimila-
49 tion and its various methodologies [16, 17] have reached a widespread and
50 worldwide interests in many federal research institutes and universities such
51 as the National Center for Atmospheric Research (NCAR, US); the National
52 Centers for Environmental Prediction (NCEP, US); the Deutscher Wetterdi-
53 enst (DWD, Germany); the Met Office (University of Reading and Imperial
54 College of London, UK); the Japan Meteorological Agency (JMA, Japan);
55 the Canadian Association of Management Consultants (CMC, Canada) and
56 the Euro- Mediterranean Center for Climate Changes (CMCC, Italy). Since
57 10 years, the Variational DA (VarDA) approaches [18, 19] have gained accep-
58 tance for its accuracy and efficiency and thus as a powerful method. VarDA
59 methodology is based on the minimisation of a function which estimates the
60 discrepancy between numerical results and observations assuming that the
61 two sources of information, forecast and observations, have errors that are
62 adequately described by error covariance matrices.

63 A POD-EnVar DA method to identify pollutant source location and wind pa-

64 rameters from observations of the gas concentration is described in [20]. The
65 POD-EnVar DA is coupled with a CFD software based on a Lattice Boltz-
66 mann Method (LBM) code and V-LES algorithm (PowerFLOW). A sensor
67 placement procedure based on global sensitivity analysis techniques has also
68 been proposed to improve the performances of the assimilation process. Us-
69 ing appropriate sensor placement, the position of the source can be identified
70 with an accuracy of only a few meters. An EnKF method to improve the pre-
71 diction of air flow (using the OpenFoam libraries as CFD software) in a real
72 urban environment using wind sensors located in Stanford’s campus, US, is
73 proposed in [21]. The location as well as the number of sensors are discussed,
74 highlighting that sensors located at roof height allows a better prediction of
75 the velocity field. Moreover, with careful selection of the sensor location,
76 their method is also able to accurately retrieve the probability distribution of
77 the inlet wind velocity and direction. Finally, an Optimal Three Dimensional
78 Variational (3DVar) data assimilation model coupled with a mesh-adaptivity
79 open-source CFD software (Fluidity) is developed in [22, 23]. The method
80 and its parametrisation is fully described and then successfully applied to
81 a real urban environment located in London, showing that the error in the
82 pollutant dispersion and the flow field can be reduced up to one order of
83 magnitude compared to before the VarDA process. Moreover, this reduction
84 in error propagates, as expected, in the next time step of the forecasted model
85 (Fluidity).

86 As mentioned before, the validation of CFD models (for urban environment

87 simulations) are usually performed by comparing, more or less successfully,
88 results to wind tunnel experiments, with a trend of higher discrepancy when
89 increasing the complexity of the urban layout. Before going towards a com-
90 parison with full-scale experiment, the coupling of DA and CFD has also
91 to be considered as a way to improve the comparison between wind tun-
92 nel and CFD results. The assimilation of pressure coefficient from a wind
93 tunnel experiment in open-source CFD software (SU2) is proposed in [24]
94 for the well-known 2D NACA 0012 and RAE 2822 airfoils. The sensitivity
95 of the results depending on the number of observation points is discussed,
96 highlighting that the assimilation works even with a very limited number of
97 measurements (4% of the original data set was used). An Ensemble Kalman
98 Filter (EnKF) method is used in [25] to assimilate values of surface pressure
99 provided by a wind tunnel experiment around a so-called “squared cylinder”
100 which can be assimilated to a single isolated building. They highlight that
101 such a coupling method is promising, however, the 3D effect of the flow is ne-
102 glected and only a 2D simulation is considered. Finally, the coupling between
103 a Monte Carlo dispersion model (probabilistic model) and an EnKF method
104 is developed in [26] showing that the error in the calculated concentration is
105 reduced when coupling with DA.

106 The work presented in this paper aims to address the discrepancy between
107 CFD results and wind tunnel data in terms of pollutant concentration pre-
108 diction in a real urban environment. Thereby, CFD will be coupled with
109 a novel data assimilation approach to show how data assimilation enhances

110 predictions and reduces the errors between measurements and simulations.
111 In this paper, the Optimal Three Dimensional Variational (3DVar) data
112 assimilation model presented in [22], which has been developed and im-
113 plemented for improving air pollution prediction, is used. The forecasted
114 model to be improved is the open-source CFD software Fluidity ([http:](http://fluidityproject.github.io/)
115 [//fluidityproject.github.io/](http://fluidityproject.github.io/)) [27], and the observed data are concen-
116 tration values from a wind tunnel experiment performed in the EnFlo Mete-
117 orological Wind Tunnel [12].

118 The CFD Large Eddy Simulation (LES) method and the Optimal 3DVar DA
119 model are first described in Section 2. The case set-up (wind tunnel exper-
120 iment and CFD simulation) is then detailed in Section 3. The results using
121 DA to improve the prediction of the pollutant concentration are presented
122 in Section 4. Finally, conclusions are provided in Section 5.

123 **2. Methodology**

124 *2.1. The Large Eddy Simulation method and Mesh Adaptivity*

125 Over the last two decades, the Large Eddy Simulation (LES) method has
126 become one of the most popular tool for atmospheric sciences, as it enables
127 a more accurate capturing of the turbulent flows compared to the traditional
128 Reynolds-Averaged Navier-Stokes (RANS) approach [10, 28, 29, 30, 31, 9].
129 The LES approach, although still complex and computationally demanding
130 is “favoured” because it allows a more accurate representation of turbulence:

131 it achieves this by separating the flow into resolved and unresolved scales
 132 based on a cut-off length scale Δ . For scales greater than Δ , the flow is
 133 resolved and numerically solved, whilst for scales smaller than Δ , the flow
 134 is unresolved and represented by a sub-grid scale model. The subgrid scale
 135 model is crucial in representing the flow of turbulent energy from the large-
 136 scale (resolved) scale motions to the smallest (unresolved) scales where energy
 137 is dissipated [32]. The importance of the subgrid scale model was very clearly
 138 noted and considered in the very early works of the development of the
 139 LES methodology - especially the need to address anisotropic filtering and
 140 inhomogeneous effects [33, 34].

141 The LES equations describing turbulent flows are based on the filtered three-
 142 dimensional incompressible Navier-Stokes (NS) equations: continuity of mass
 143 (equation (1)) and momentum equations (equation (2)) [35].

$$\nabla \cdot \bar{u} = 0 \tag{1}$$

$$\frac{\partial \bar{u}}{\partial t} + \bar{u} \cdot \nabla \bar{u} = -\frac{1}{\rho} \nabla \bar{p} + \nabla \cdot [(\nu + \nu_\tau) \nabla \bar{u}] \tag{2}$$

144 where \bar{u} is the resolved velocity (m/s), \bar{p} is the resolved pressure (Pa), ρ is
 145 the fluid density (kg/m³), ν is the kinematic viscosity (m²/s) and ν_τ is the
 146 anisotropic eddy viscosity (m²/s).

147 The subgrid-scale model in Fluidity is based on the Smagorinsky model in
 148 which the eddy viscosity ν_τ is expressed by equation (3).

$$\nu_\tau = C_S^2 \Delta^2 |\overline{S}| \quad (3)$$

149 C_S is the Smagorinsky coefficient (taken equal to 0.1), Δ is the Smagorinsky
 150 length scale which depends on the local element size and $|\overline{S}|$ is the strain rate
 151 expressed as in equation (4).

$$|\overline{S}| = (2\overline{S}_{ij}\overline{S}_{ij})^{1/2} \quad (4)$$

152 where \overline{S}_{ij} is the local strain rate defined by equation (5).

$$\overline{S}_{ij} = \frac{1}{2} \left(\frac{\partial \overline{u}_i}{\partial x_j} + \frac{\partial \overline{u}_j}{\partial x_i} \right) \quad (5)$$

153 A novel component in the implementation of the subgrid-scale model within
 154 Fluidity is the anisotropic eddy viscosity tensor [35] defined as in equation (6):

$$\nu_\tau = 4C_S^2 |\overline{S}| \mathcal{M}^{-1} \quad (6)$$

155 where \mathcal{M} is the length scale metric from the adaptivity process [36] used
 156 here to relate eddy viscosity to the local grid size as shown in equation 7.

$$\mathcal{M}^{-1} = V^T \begin{pmatrix} h_\zeta^2 & 0 & 0 \\ 0 & h_\eta^2 & 0 \\ 0 & 0 & h_\xi^2 \end{pmatrix} V \quad (7)$$

157 with V^T and V the rotational transformations to transform from the local
 158 to the global coordinate systems and (h_ζ, h_η, h_ξ) the local element sizes. The
 159 factor of 4 arises because the filter width separating resolved and unresolved
 160 scales is assumed to be twice the local element size, which is squared in the
 161 viscosity model. It has been shown that an anisotropic eddy viscosity gives
 162 better results for flow simulations on unstructured grids [35].

163 The transport of a scalar field c (i.e, a passive tracer) in kg/m^3 is expressed
 164 using a classic advection-diffusion equation having a source term as in equa-
 165 tion (8):

$$\frac{\partial c}{\partial t} + \nabla \cdot (\mathbf{u}c) = \nabla \cdot (\bar{\kappa} \nabla c) + F \quad (8)$$

166 where \mathbf{u} is the velocity vector (m/s), $\bar{\kappa}$ is the diffusivity tensor (m^2/s) and
 167 F represents the source terms ($\text{kg}/\text{m}^3/\text{s}$).

168 The source term F is expressed by equation (9):

$$F = \frac{Q\rho}{V} \quad (9)$$

169 where Q is a volumetric flow rate expressed in m^3/s and V is the volume of
170 the source in m^3 .

171 The behaviour of the atmospheric boundary layer in Fluidity is represented
172 using a turbulent inlet velocity based on a synthetic eddy method [37, 38].
173 The turbulent inlet velocity is controlled by: the turbulence length scales pro-
174 files (L_u, L_v, L_w) , the mean velocity profiles $(\bar{u}, \bar{v}, \bar{w})$ as well as the Reynolds
175 stresses profiles $(\overline{u'u'}, \overline{v'v'}, \overline{w'w'})$.

176 The need for combining the LES approach with adaptive meshes has been
177 tackled as a way of overcoming the large range of length scales that exist in
178 turbulent flows [39]. The challenge of combining the LES approach with 3D,
179 adaptive, unstructured meshes was first undertaken and implemented within
180 the Fluidity software [36, 35, 27]. Hence, one of the key and innovative as-
181 pects of Fluidity is its mesh-adaptivity capability on unstructured meshes.
182 The adaptivity process allows: (i) the addition or reduction of the number
183 of nodes and elements, leading subsequently to refining or coarsening of the
184 mesh depending on the area of interest; (ii) smoothing of the mesh by mov-
185 ing nodes whilst keeping the overall number of elements and nodes the same.
186 A-posteriori error estimates are made, aiming at certain targets for error [27].
187 Adaptivity options can be field-specific, i.e. different computed fields can be
188 configured with their own specific adaptivity options. This process allows
189 to have fine mesh in region where small-scale and important physical pro-
190 cesses occur, while keeping a coarser mesh elsewhere, and then allowing to

191 considerably reduce the total computation time [36].

192 2.2. Data Assimilation

193 Let n be a fixed time level and let c^n be the state variable c as described in
194 equation (8) at the fixed time level n . Let v^n be an observation of the state
195 variable at time n and let consider a mapping H as in equation 10.

$$H : c^n \mapsto v^n. \quad (10)$$

196 Let $d^n = v^n - H(c^n)$ be the misfit. In this section, we introduce a Data
197 Assimilation process in which the solution of the forecasting model (Fluidity)
198 obtained from equation (8) is combined with information provided by a wind
199 tunnel experiment in order to improve the accuracy of the solution c^n , i.e. to
200 reduce d^n . The aim of the Data Assimilation problem is to find an optimal
201 trade-off between the prediction made based on the Fluidity system state
202 c^n (background) defined in equation (8) and the available observation v^n
203 provided by the wind tunnel.

204 For a fixed time step n , given c^n and v^n , the DA process consists in finding
205 c^{DA} as an inverse solution of equation (11) subject to the constraint given
206 by equation (12).

$$v^n = H(c^{DA}), \quad (11)$$

$$c^{DA} = c^n. \quad (12)$$

207 Since H is typically rank deficient, the equation (11) is an ill-posed inverse
 208 problem [40, 41]. The Tikhonov formulation [42] leads to an unconstrained
 209 least squares problem, where the term in equation (12) provided by Fluidity
 210 ensures the existence of a unique solution of equation (11). The DA process
 211 can be then described as following [43]:

$$c^{DA} = \operatorname{argmin}_c \{ \|c - c^n\|_{\mathbf{B}^{-1}}^2 + \|v^n - H(c)\|_{\mathbf{R}^{-1}}^2 \} \quad (13)$$

212 where \mathbf{R} and \mathbf{B} are the observation and model error covariance matrices
 213 respectively defined by equation (14) and equation (15):

$$\mathbf{R} := \sigma_0^2 \mathbf{I} \quad (14)$$

214 with $0 \leq \sigma_0^2 \leq 1$ representing the variance value of the distribution of the
 215 instruments errors and \mathbf{I} the identical matrix;

$$\mathbf{B} = \mathbf{V}\mathbf{V}^T \quad (15)$$

216 where \mathbf{B} is the background error covariance matrix as defined in Definition 1
 217 associated with the state c since the true state will differ from the simulated

218 state by random or systematic errors.

219 **Definition 1** (Variance-Covariance Matrix). *Let \mathbf{X} be a matrix of measure-*
 220 *ments of pv physical variables at spatial locations $\mathcal{D} = \{x_j\}_{j=1,\dots,np}$ for a*
 221 *correlation time window $[0, T_1] = \{\tau_k\}_{k=1,\dots,M}$:*

$$\mathbf{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_{NP} \end{bmatrix} \in \mathfrak{R}^{NP \times M} \quad (16)$$

222 *where each of NP row is a time series for a given location and $NP = [pv] \cdot np$.*
 223 *Let's assume that each row X_i of \mathbf{X} has a mean $E[X_i] = \{m_i\}_{i=1,\dots,NP}$ and*
 224 *let's define $\mathbf{m} = (m_i)_{i=1,\dots,NP}$. Hence, the deviation matrix is:*

$$\mathbf{V} = \mathbf{X} - \mathbf{m} \in \mathfrak{R}^{NP \times M}, \quad (17)$$

225 *If each vector X_i has a distribution with probability density function P , then*
 226 *the expected value of X_i is defined by:*

$$E(X_i) = \frac{1}{M-1} \sum_{j=1,\dots,M} x_{ij} P(X_j) \quad . \quad (18)$$

227 *The variance-covariance matrix $\mathbf{B} \in \mathfrak{R}^{NP \times NP}$ of \mathbf{X} (equation (19)) is then*
 228 *defined via the expected value of the outer product:*

$$\mathbf{B} = \mathbf{V}\mathbf{V}^T \quad . \quad (19)$$

229 If equation (13) is linearised around the background state [44], it yields:

$$c = c^n + \delta c \quad (20)$$

230 where $\delta c = c - c^n$ denotes the increments. The DA problem can then be
 231 re-formulated by the following form:

$$\delta c^{DA} = \underset{\delta c}{\operatorname{argmin}} \left\{ \frac{1}{2} \delta c^T \mathbf{B}^{-1} \delta c + \frac{1}{2} (\mathbf{H} \delta c - d^n)^T \mathbf{R}^{-1} (\mathbf{H} \delta c - d^n) \right\} \quad (21)$$

232 where

$$d^n = v^n - H(c^n) \quad (22)$$

233 is the misfit between the observation and the solution computed by Fluidity
 234 and

$$H(c) \simeq H(c^n) + \mathbf{H} \delta c \quad (23)$$

235 denotes the linearised observational and model operators evaluated at $c = c^n$
 236 where \mathbf{H} is the Hessian of H .

237 In equation (21), the minimisation problem is defined on the field of incre-
 238 ments [45]. In order to avoid the inversion of \mathbf{B} , as $\mathbf{B} = \mathbf{V}\mathbf{V}^T$ (see equa-
 239 tion (19)), the minimisation can be computed with respect to a new variable
 240 $w = \mathbf{V}^+ \delta c$ [44], where \mathbf{V}^+ denotes the generalised inverse of \mathbf{V} , yielding to:

$$w^{DA} = \underset{w}{\operatorname{argmin}} \left\{ \frac{1}{2} w^T w + \frac{1}{2} (\mathbf{H}\mathbf{V}w - d^n)^T \mathbf{R}^{-1} (\mathbf{H}\mathbf{V}w - d^n) \right\} \quad (24)$$

241 As the background error covariance matrix is ill-conditioned [41], in order
 242 to improve the conditioning, only Empirical Orthogonal Functions (EOFs)
 243 of the first largest eigenvalues of the error covariance matrix are considered.
 244 Since its introduction to meteorology by Edward Lorenz [46], EOFs analysis
 245 has become a fundamental tool in atmosphere, ocean, and climate science for
 246 data diagnostics and dynamical mode reduction. Each of these applications
 247 exploits the fact that EOFs allow a decomposition of a data function into a
 248 set of orthogonal functions, which are designed so that only a few of these
 249 functions are needed in lower-dimensional approximations [47]. Furthermore,
 250 since EOFs are the eigenvectors of the error covariance matrix [48], its con-
 251 dition number is reduced as well. Nevertheless, the accuracy of the solution
 252 obtained by truncating EOFs exhibits a severe sensibility to the variation
 253 of the value of the truncation parameter, so that a suitably choice of the
 254 number of EOFs is strongly recommended. This issue introduces a severe
 255 drawback to the reliability of EOFs truncation, hence to the usability of the
 256 operative software in different scenarios [48, 49]. In this paper, we set the
 257 optimal choice of the truncation parameter as a trade-off between efficiency
 258 and accuracy of the DA algorithm as introduced in [22].

259 The Optimal 3DVar data assimilation model as implemented in this paper is
 260 summarised in Algorithm 1.

Algorithm 1 : A VarDA algorithm to assimilate Wind Tunnel data into Fluidity.

- 1: Input: $\alpha, \{v_k\}_{k=0,\dots,m}, c_0^M$
 - 2: Define \mathbf{H}
 - 3: Compute $d_k \leftarrow v - \mathbf{H}c_0^M$ ▷ compute the misfit
 - 4: Define \mathbf{R} starting from the wind tunnel data v
 - 5: Compute \mathbf{V} ▷ deviance matrix defined in (17)
 - 6: Compute $\mathbf{V}_\tau = EOFs(\mathbf{V}, \tau)$ ▷ reduced space computed by EOFs
 - 7: Define the initial value of $\delta\mathbf{u}^{DA}$
 - 8: Compute $w \leftarrow \mathbf{V}_\tau^+ \delta c^{DA}$ ▷ from the physical to reduced space
 - 9: repeat ▷ start of the L-BFGS steps
 - 10: Compute $J \leftarrow J(w)$
 - 11: Compute $gradJ \leftarrow \nabla J(w)$
 - 12: Compute new values for w
 - 13: until (Convergence on w is obtained) ▷ end of the L-BFGS steps
 - 14: Compute $\delta c^{DA} \leftarrow \mathbf{V}_\tau w$ ▷ from the reduced to physical space
 - 15: Compute $c^{DA} \leftarrow c_0^M + \delta c^{DA}$
-

261 **3. Case Set-up**

262 Initial validations of Fluidity have already been reported in which compar-
 263 isons of velocity, mean pollutant concentration predictions and surface pres-
 264 sures with wind tunnel data were carried out [11, 12, 13, 15, 35]. How-
 265 ever, comparisons between experiments and simulations are most of the time
 266 plagued by discrepancies. In [12], the comparison of mean pollutant concen-
 267 trations at 81 detector locations was carried out and it was observed that
 268 the errors between simulations and measurements ranged between 3% to over
 269 50%. Thereby, the same test case than in [12] is used in this paper and is
 270 coupled with the Optimal Three Dimensional Variational (3DVar) data as-
 271 simulation model presented in Section 2.2 (and fully described in [22]) in order

Building	Height (cm)
N	14.28
1	13.15
2	12.38
3	11.52
4	3.15
5	9.71
7	12.28

Table 1: Buildings heights, used in the LES simulation, based on the wind tunnel configuration. The buildings labels refer to the ones given in Figure 1b.

272 to improve the accuracy of the results predicted by Fluidity.

273 3.1. Geometry

274 A 7-buildings configuration is considered in this paper as shown in Figure 1.

275 The buildings represent a real, small neighbourhood area in central London,

276 UK (51°30'00.0"N, 0°12'00.9"W), at the scale of 1:200 (wind tunnel scale).

277 The heights of the seven buildings are given in Table 1, where the labels of

278 each building refer to the ones given in Figure 1b.

279 3.2. Wind tunnel data

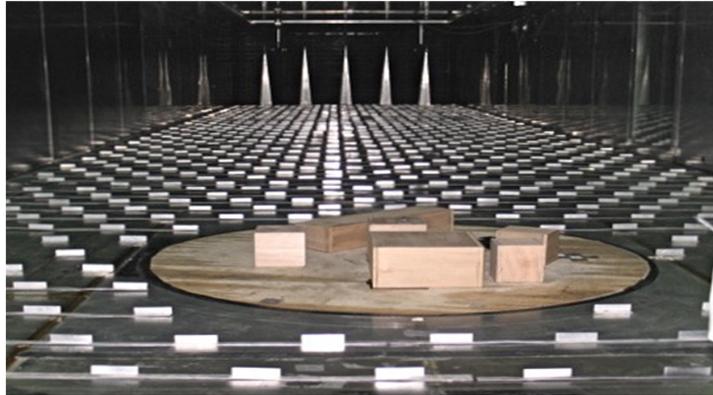
280 A set of experiments were carried out at the EnFlo wind tunnel [12] for

281 the 7-buildings configuration (Figure 1a). The geometry represented is at

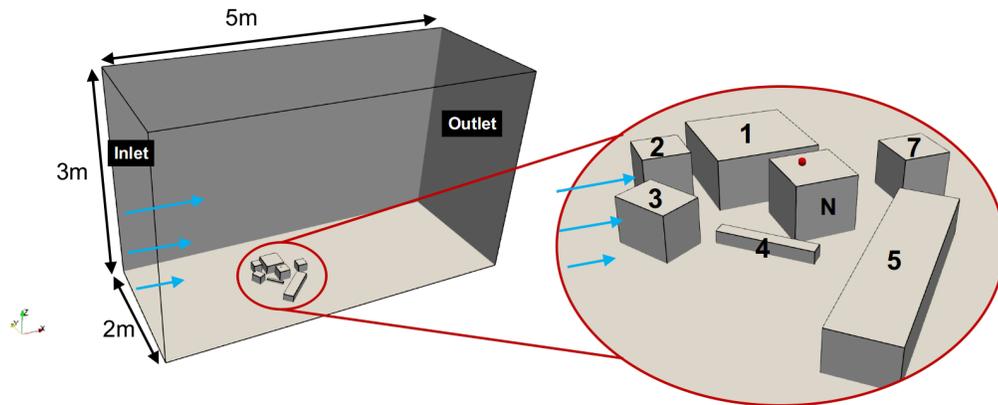
282 1:200 scale. The experiments were carried out in a fully developed, 1m-deep,

283 simulated atmospheric boundary layer with a reference wind velocity U_{ref}

284 of 2 m/s. The experimental atmospheric boundary layer represents neutral



(a) Wind tunnel set-up



(b) CFD Computational Domain

Figure 1: The 7-buildings configuration (a) in the wind tunnel experiment and (b) in the CFD simulation. In (b) the location of the source is denoted by the red sphere at the top of Building N and the wind direction is shown by the blue arrows.

285 atmospheric conditions and is initiated by a set of Irwin spires (vorticity-
286 generators) at the inlet of the wind tunnel working section, with roughness
287 elements on the floor to maintain the surface roughness condition. The sur-
288 face roughness length z_0 and the friction velocity u^* are equal to 1.5 mm
289 and 0.057 m/s, respectively, with u^* being the air velocity at the edge of the
290 boundary layer.

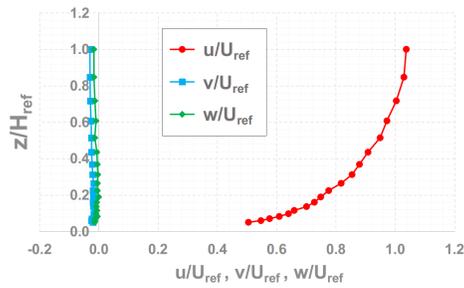
291 A passive tracer (propane) was emitted from a horizontal source, having a
292 diameter of 20 mm, above Building N (Figure 1b) at 15.08 cm height from
293 the ground of the test section, i.e 0.8 cm above Building N (having an height
294 of 14.28 cm). It has to be noted that the source is not centred on the
295 top of Building N. The tracer release flow rate in experiments was equal to
296 $Q_{WT} = 2.4 \times 10^{-7} m^3/s$. The assumption could be made that there is no den-
297 sity difference between the emission gas and the surrounding fluid (air) [50].
298 Indeed, the propane gas (the trace gas) is diluted into the surrounding air
299 such that the percentage proportion of propane/trace gas ranges between
300 0.99% to 2.1% of the total released gas. This mixture is considered neu-
301 trally buoyant and is released at a point source. These proportions and
302 this gas are commonly used in wind tunnel experiment as non-reactive and
303 non-depositing tracer gas, so that it disperses as a passive tracer in the
304 flow [50, 13]. Due to the large amount of air mass, it is considered that the
305 trace particle number is small so that the trace particles do not significantly
306 influence the density. The density of the emission is then considered to be
307 the same as of the surrounding air.

308 Mean tracer concentrations were measured using Combustion Fast Flame
309 Ionisation Detectors (FFIDs) carried on a three-dimensional traverse system
310 and each point measurement is an average over an acquisition period of 2
311 minutes. Measurements were taken for varying wind directions and model
312 configuration, however only one configuration and one wind direction is con-
313 sidered in this paper. The tracer concentration was obtained at 738 different
314 locations, located downstream the source.

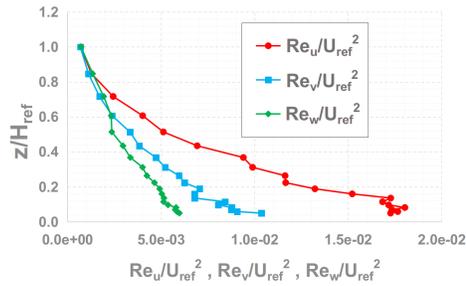
315 Figure 2 shows the inflow profiles of the velocity, the Reynolds stresses and
316 the turbulence length scales used in the wind tunnel experiment, where the
317 reference height H_{ref} is the boundary layer height and U_{ref} the reference
318 velocity. The Reynolds number based on the mean building height $H_{mean} =$
319 $10.9cm$ is approximately equal to 1.4×10^4 .

320 3.3. The LES Computational set-up

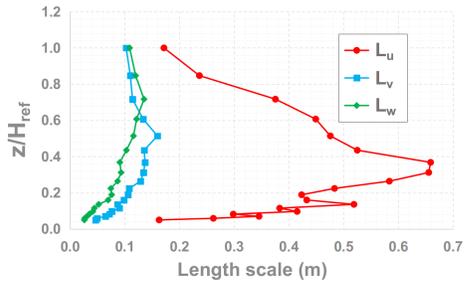
321 The Fluidity-LES software was used to model the flow field and the passive
322 tracer concentrations within the 7-buildings configuration. The dimensions of
323 the computational domain covered a volume of $5.0\text{ m} \times 2.0\text{ m} \times 3.0\text{ m}$ ($x \times y \times z$
324 -direction) as shown in Figure 1b, allowing a relatively long-development
325 section for the formation of a deep boundary layer in the LES simulation.
326 The blockage ratio is equal to 2.3%, below the maximum value recommended
327 of 3% [51, 52]. The height, the width and the length of the domain are more
328 than 5 times higher than the taller building (Building N) and/or the diameter



(a) Velocity profiles



(b) Reynolds stresses profiles



(c) Length scale profiles

Figure 2: Inflow profiles of the three components of (a) the velocity, (b) the Reynolds stresses and (c) the turbulence length scales used in both the wind tunnel experiments and the CFD simulations. H_{ref} is the boundary layer depth and U_{ref} is the reference velocity.

329 of the buildings area, hence following the guidance rules for CFD in urban
330 environment [51, 52].

331 The location of the inlet and the outlet of the domain are shown in Fig-
332 ure 1b. A turbulent velocity is prescribed at the inlet, based on a synthetic
333 eddy method [37] and the blue arrows in Figure 1b shows the wind direction.
334 The wind direction is directly perpendicular to the front façades of Buildings
335 1, 2 and N. The mean velocities, the turbulence length scales as well as the
336 Reynolds stresses profiles are set-up using the profiles provided by the wind
337 tunnel experiments as shown in Figure 2. In a real urban dispersion prob-
338 lem, wind direction and velocity are constantly changing, however this is not
339 taken into account here as the application is proposed for wind tunnel test
340 cases only. Indeed, in wind tunnel, the experiments are done in controlled
341 environments where the wind direction and velocity are fixed. The down-
342 stream boundary (outlet) is left as pressure boundary, whilst the remaining
343 boundary conditions consisted of: (i) the “no slip” condition for the solid
344 walls of buildings and “floor” of the domain, and (ii) the “slip/no shear”
345 condition for the free surfaces (sides and top of the domain).

346 The emission source was placed at the top of the central building, i.e Building
347 N, at the same location and height than in the wind tunnel as shown by
348 the red sphere in Figure 1b. The diameter of the source is equal to 20
349 mm and the diffusion coefficient of propane in an excess of air is set to
350 $1 \times 10^{-5} m^2/s$. **The propane is considered as non-reactive and non-depositing**

351 tracer gas [50]. Thus, the propane emission in the simulations is considered
 352 as a passive tracer, i.e. no density effect/variation with the surrounding fluid
 353 and travel with the air flow velocity such that the classic advection-diffusion
 354 with a source term (equation 8) is used. The source term F , expressed by
 355 equation (9), is set equal to $F = 1kg/m^3/s$, leading to a volumetric flow
 356 rate Q_F equal to $2.5 \times 10^{-6}m^3/s$, i.e. one order of magnitude higher than
 357 in experiments. In order to be compared, the concentration c from wind
 358 tunnel experiment and the ones obtained from Fluidity, the concentrations
 359 are commonly normalised using equation (25) [53, 54]:

$$c^* = \frac{cU_{ref}H_{mean}^2}{Q} \quad (25)$$

360 where c^* is the normalised concentration, U_{ref} is the reference velocity (m/s)
 361 at the top of the boundary layer and H_{mean} is the mean building height.
 362 U_{ref} and H_{mean} are the same in both experiment and simulation. Hence, the
 363 concentrations from the wind tunnel c_{WT} are converted into their equivalent
 364 “Fluidity” values v^n using equation (26):

$$v^n = c_{WT} \frac{Q_F}{Q_{WT}} \quad (26)$$

365 where Q_F and Q_{WT} stand for the volumetric flow rate in Fluidity and in
 366 wind tunnel, respectively. In the DA process, this modified wind tunnel
 367 concentration corresponds to the observed data, i.e v^n .

368 All the equations are solved using second order schemes in time and space.
369 The NS equations are discretised using a continuous Galerkin discretisa-
370 tion, while the advection-diffusion is discretised using a second order upwind
371 scheme. An adaptive time step is used and the CFL number is equal to 0.9,
372 leading to an average time step equal to 1×10^{-3} , while the Crank-Nicholson
373 scheme is used for the time discretisation. Note that the time step is not
374 constant in the simulation because of the use of mesh adaptivity. Absolute
375 and relative convergence errors were set to 10^{-12} and 10^{-7} , respectively for
376 all fields (pressure, velocity and tracer).

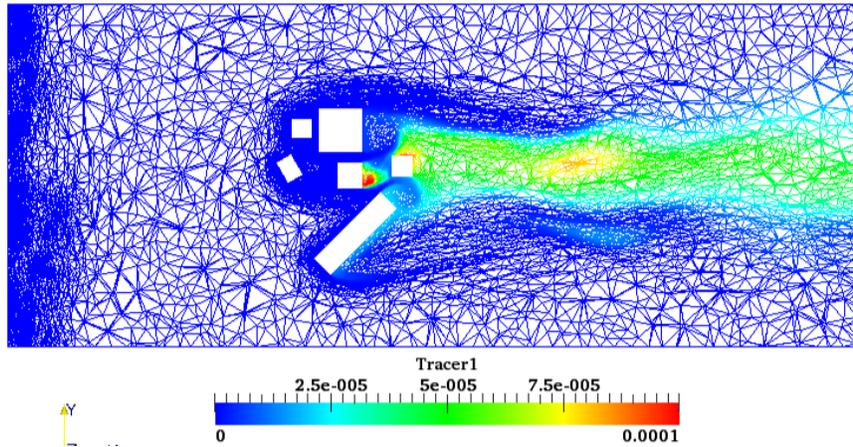
377 *3.4. Mesh adaptivity and Supermesh*

378 For the LES simulation presented in this work, field-specific adaptivity op-
379 tions were assigned to the velocity field and the tracer field. For both fields,
380 mesh resolution was also controlled by specifying the maximum and the min-
381 imum sizes of the elements in the domain. They are respectively taken equal
382 to 1 cm and 15 cm. Moreover, to resolve the source, the mesh is locally
383 controlled around the source location by setting the minimum edge length to
384 be 3 mm, and allowing the maximum element size to be 4 mm. The mesh
385 was adapted every 15 time steps, and anisotropic gradation was also allowed
386 in the simulation. The maximum number of nodes was set to 400,000. An
387 example of the adaptivity effect on the computational mesh can be seen in
388 Figure 3 for the instantaneous tracer field on two horizontal planes. The res-
389 olution of the mesh is fine near the inlet to capture accurately all the eddies

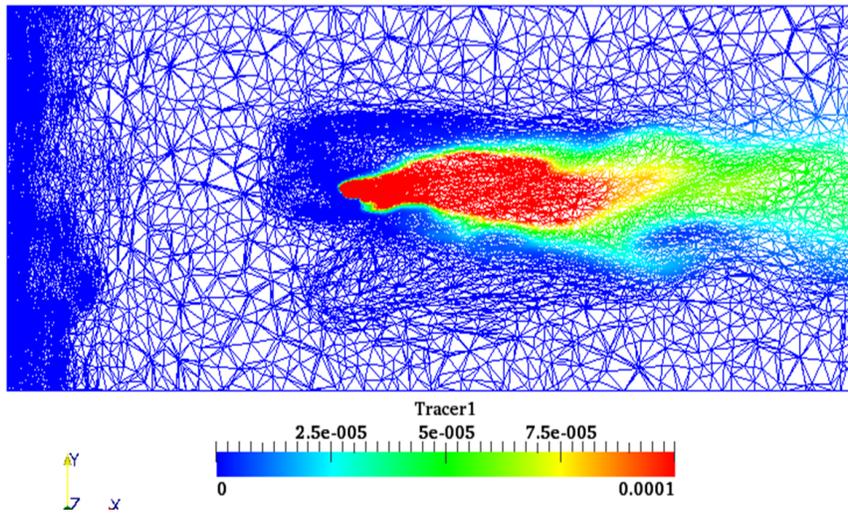
390 coming into the domain: this is a direct effect of the mesh adaptivity. It has
391 to be noted that the mesh is changing every 15 time steps, and the meshes
392 shown in Figure 3 is an example of instantaneous mesh.

393 In order to compute \mathbf{B} , the background error covariance matrix, (equa-
394 tion (15)), the number of nodes in the mesh has to be “fixed”, i.e always
395 the same at every time step. Thereby, Fluidity is running with mesh adap-
396 tivity for 34 sec (real time), which is sufficient for the flow statistics to reach
397 a quasi-steady state. From this point onwards, the mesh obtained has a min-
398 imum and maximum edge lengths equal to 1.4 mm and 25 cm, respectively;
399 while the number of nodes in the mesh is equal to 170,775. This mesh will
400 be referred as the *supermesh* in the following and is shown in Figure 3. The
401 *supermesh* is considered as an optimal mesh, as the simulation has run long
402 enough to have fine elements in areas where important physical processes
403 occur repeatedly. Fluidity results are then projected onto that *supermesh* in
404 order to compute \mathbf{B} , the background error covariance matrix, (equation (15)).
405 It has to be mentioned that this process (projection of all data) has to be
406 done only once, as \mathbf{B} has to be computed only once. The mesh adaptivity
407 process can then be used normally when Fluidity is running: the projection
408 of Fluidity data onto the *supermesh* is then done only for the time step at
409 which observed data want to be assimilated.

410 The wind tunnel data v^n , i.e the observed data, has also to be projected on
411 the *supermesh*. The location of sensors in wind tunnel does not necessarily



(a) $z = 6.5$ cm



(b) $z = 14.8$ cm

Figure 3: Instantaneous tracer concentration, i.e. pollutant concentration, at $t = 34$ sec for horizontal planes (xOy) at heights (a) $z = 6.5$ cm and (b) $z = 14.8$ cm, obtained from Fluidity. The mesh shown in these figures also corresponds to the mesh used as the *supermesh*. The tracer concentration ranges between $0\text{kg}/\text{m}^3$ (blue colour) and $1 \times 10^{-4}\text{kg}/\text{m}^3$ (red colour).

412 lie on a *supermesh* node. Therefore, using interpolation method, the sensor
413 value is distributed to the four nodes of the tetrahedron in which lies the
414 sensor. As one mesh node can be part of several tetrahedrons in which lie
415 different sensors, the number of nodes in the mesh where sensors data are
416 assigned is smaller than four times the number of sensors. This process leads
417 to a number of nodes equal to 1391, i.e. values from wind tunnel experiments
418 are assigned to 1391 nodes in the *supermesh*.

419 4. Results and Discussion

420 A comparison between Fluidity results and wind tunnel data for 81 detec-
421 tors was carried out in [12], with differences/errors between simulations and
422 measurements ranging between 3% to over 50%. The results presented here
423 aim to reduce these errors using the DA method described in Section 2.2. In
424 this section, 1391 observation points, located downstream of the pollutant
425 source, are considered and their locations are shown in Figure 4.

426 4.1. Accuracy evaluation

427 The accuracy of the DA results are evaluated using:

- 428 • the absolute error

$$E(c) = |c - v^n| \tag{27}$$

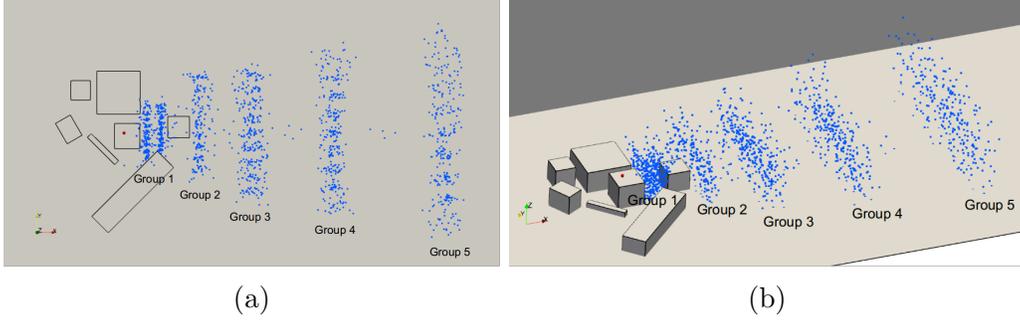


Figure 4: Location (blue dots) of the sensors in the domain. Five groups of sensors can be identified, based on their distances from the source. The red sphere denotes the location of the source.

- 429 • the mean squared error

$$MSE(c) = \frac{\|c - v^n\|_{L^2}}{\|v^n\|_{L^2}} \quad (28)$$

430 where c is either c^n the Fluidity concentration at time step n or c^{DA} the
 431 corrected concentration using DA (Algorithm 1) and v^n is the wind tunnel
 432 observed data.

433 Figure 5 shows the values of the absolute errors $E(c^n)$ and $E(c^{DA})$ on three
 434 different slices: through the oriented planes (xOy) , (xOz) and (yOz) . After
 435 the DA process, the absolute error is visibly reduced by almost one order
 436 of magnitude everywhere in the domain. The absolute error $E(c^n)$ ranges
 437 between 1×10^{-5} and 3×10^{-6} , with error values decreasing as the distance in
 438 the y -direction from the source increases (Figure 5e). After the DA process,
 439 the absolute error $E(c^{DA})$ becomes lower than 2.5×10^{-6} at every sensor
 440 location.

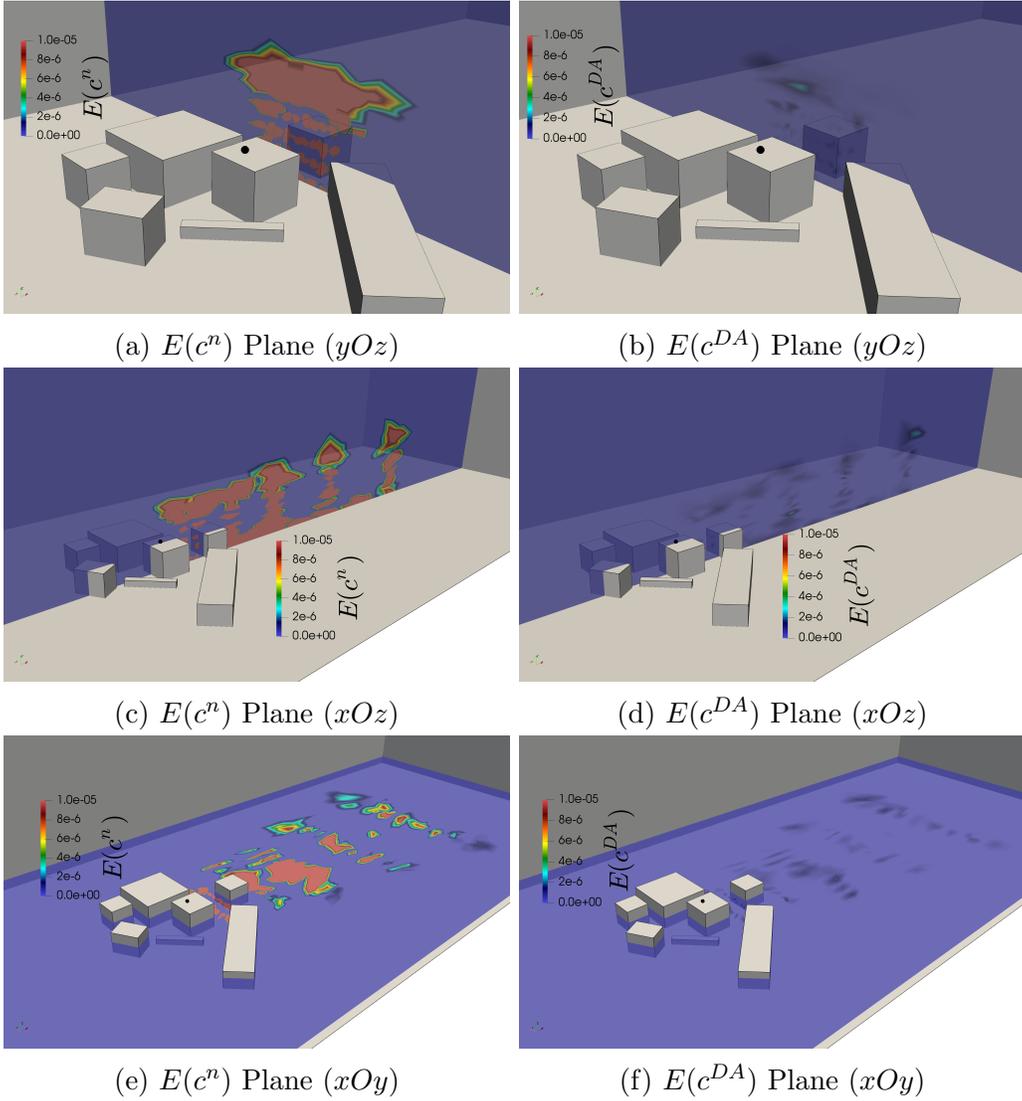


Figure 5: Values of the absolute error $E(c)$ (equation (27)) through three slices: through a plane (xOy), a plane (xOz) and a plane (yOz). The absolute errors shown are computed before ($E(c^n)$) and after ($E(c^{DA})$) the assimilation process. The scale of $E(c)$ ranges between 0 (blue colour) and 1×10^{-5} (red colour) in all sub-figures. The black sphere denotes the source location.

441 Figure 6a shows the variation of the mean squared errors $MSE(c^n)$ and
 442 $MSE(c^{DA})$ as a function of the number of assimilated observations. The
 443 $MSE(c^n)$ does not depend on the number of observations and is then con-
 444 stant and equal to 3.749. $MSE(c^{DA})$ decreases as a function of the number
 445 of observations, reaching a value of 5.6×10^{-3} for 1391 observations assimi-
 446 lated: the DA process allows a reduction of the mean squared error by almost
 447 three order of magnitudes. $MSE(c^{DA})$ is reduced by one order of magnitude
 448 (3.75×10^{-1}) and two order of magnitudes (3.75×10^{-2}) assimilating 722
 449 and 1312 observations, respectively. Moreover, $MSE(c^{DA})$ is approximately
 450 divided by two for 164 observations assimilated. Indeed, as shown in Fig-
 451 ure 6a, while the number of observations starts to increase, the $MSE(c^{DA})$
 452 firstly decreases very sharply, exhibiting a value of 6.7×10^{-1} for 400 number
 453 of observations. After what, the $MSE(c^{DA})$ continues to be reduced as the
 454 number of observations raises, but less quickly. The observed values v^n are
 455 assimilated in “ascending order” in terms of distance from the source, i.e from
 456 the sensor being the closer of the pollutant source to the sensor being the
 457 farthest. In other words, during the assimilation process, while the number
 458 of observations increases, more and more sensors located far way from the
 459 source are taken into account. The trend of $MSE(c^{DA})$ shown in Figure 6a
 460 tends to highlight that the closest sensors have an higher impact on the error
 461 reduction.

462 A scatter plot of the computed concentrations c^n and c^{DA} (using 1391 obser-
 463 vations) as a function of the observed data v^n is shown in Figure 6b. Ideally,

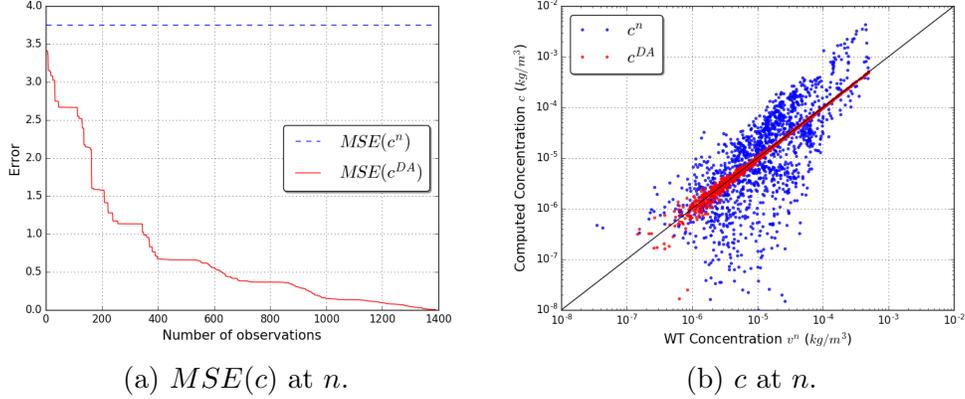


Figure 6: Values of (a) the mean squared error $MSE(c)$ (equation (28)) as a function of the number of assimilated observations and (b) the pollutant concentration c as a function of the wind tunnel data v^n (logarithm scale for both axis). In (b), the black line corresponds to the ideal matching between data and c^{DA} is obtained assimilating 1391 observations. Results at n , i.e before re-running Fluidity.

464 values should match the black line shown in Figure 6b. However, Fluidity
 465 concentration c^n are spread above and below it (blue dots in Figure 6b) with
 466 a tendency of larger spread towards low concentrations. After using the DA
 467 process, the corrected concentrations c^{DA} exhibit an obvious better agree-
 468 ment with wind tunnel data (red dots in Figure 6b), with corrected values
 469 much closer to the ideal matching (black line). The DA process performs very
 470 well in correcting Fluidity results for high concentrations, i.e concentrations
 471 higher than $1 \times 10^{-5} kg/m^3$; while, even if obviously a better agreement exists,
 472 small discrepancies still subsist for lower concentrations, i.e. concentrations
 473 lower than $1 \times 10^{-5} kg/m^3$.

474 *4.2. Impact of DA on Fluidity results*

475 In this section, the impact of DA on Fluidity results are discussed: the
476 corrected concentrations c^{DA} are used to re-run Fluidity such that these new
477 values are used as initial condition for the tracer.

478 Let $M_{n,n+i}$ denote the Fluidity software from the time step n to the time
479 step $n + i$ such that $c^{n+i} = M_{n,n+i}(c)$. The mean squared error at time step
480 $n + i$ is then defined as in equation (29):

$$MSE(M_{n,n+i}(c)) = \frac{\|M_{n,n+i}(c) - v^{n+i}\|_{L^2}}{\|v^{n+i}\|_{L^2}} \quad (29)$$

481 where c is either c^n the Fluidity concentration at time step n or c^{DA} the
482 corrected concentration using DA (Algorithm 1) and v^{n+i} is the wind tunnel
483 observed data.

484 Figure 7a shows the variation of the mean squared errors $MSE(M_{n,n+1}(c^n))$
485 and $MSE(M_{n,n+1}(c^{DA}))$ obtained after re-running Fluidity for one more time
486 step. As for $MSE(c^n)$, $MSE(M_{n,n+1}(c^n))$ does not depend of the number
487 of observations and is equal to 3.747. Figure 7a confirms that the error
488 $MSE(M_{n,n+1}(c^{DA}))$ also decreases as the number of observed data increases
489 with almost the same trend than the reduction of $MSE(c^{DA})$. For 1391
490 observations, $MSE(M_{n,n+1}(c^{DA}))$ is equal to 8.7×10^{-2} , i.e the error is re-
491 duced by two order of magnitudes compared to $MSE(M_{n,n+1}(c^n))$. It has
492 to be noted that the minimum value of $MSE(M_{n,n+1}(c^{DA}))$ is one order of

493 magnitude higher than the minimum of $MSE(c^{DA})$: this is not surprising as
494 the Fluidity software $M_{n,n+1}$ introduces intrinsically new errors. Notewor-
495 thy values can be mentioned: as for $MSE(c^{DA})$, the error is almost divided
496 by two for 164 observations assimilated ($MSE(M_{n,n+1}(c^{DA})) = 1.95$) and
497 $MSE(M_{n,n+1}(c^{DA}))$ is reduced by one order of magnitude (3.7×10^{-1}) when
498 742 observations are considered.

499 Fluidity is re-run for 200 more time steps in order to see how the reduc-
500 tion in error gained by the DA process at time step n propagates into the
501 model through time. The values of $MSE(M_{n,n+i}(c))$ as a function of the
502 time step i is shown in Figure 7b. $MSE(M_{n,n+i}(c^n))$ slightly changes over
503 time but stays however more or less constant with an average value of 3.724.
504 $MSE(M_{n,n+i}(c^{DA}))$ increases while the time step increases, tending to reach
505 the value of $MSE(M_{n,n+i}(c^n))$ after a long enough time: this behaviour is
506 expected as the model introduces new errors. This is because the physical
507 system does not change after the assimilation process as this only affects
508 the state, the boundaries and initial conditions. These are errors intrinsic
509 to the forecasting model problem which propagate on time steps. These are
510 the approximation errors introduced by the linearisation, the discretisation,
511 the model reduction... These occur when infinite-dimensional equations are
512 replaced by a finite dimensional system (that is the process of discretisa-
513 tion), or when simpler approximations to the equations are developed (e.g.,
514 by model reduction). Finally, given the numerical problem, the algorithm is
515 developed and implemented as a mathematical software. At this stage, the

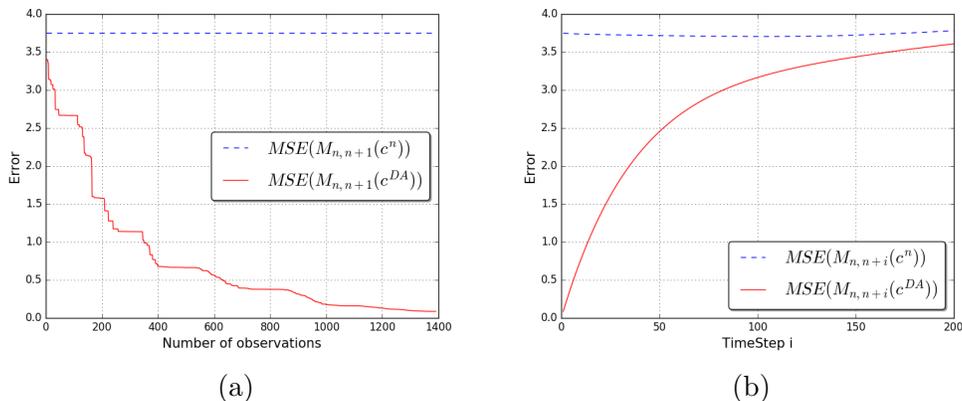


Figure 7: Values of the mean squared error $MSE(M_{n,n+i}(c))$ (equation (29)) after re-running Fluidity. (a) Variation of $MSE(M_{n,n+1}(c))$ as a function of the number of assimilated observations for $i = 1$, i.e. after one time step. (b) Variation of $MSE(M_{n,n+i}(c))$ as a function of the time step i . c^{DA} is obtained assimilating 1391 observations.

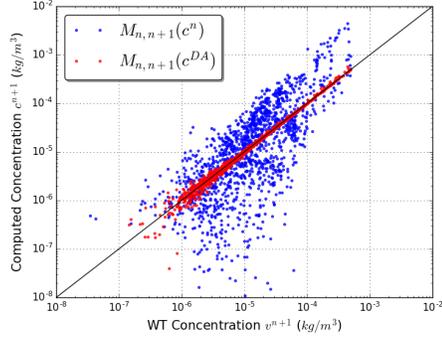
516 inevitable rounding errors introduced by working in finite-precision arith-
517 metic occurs. These errors cannot be controlled but, after few time steps,
518 the DA process can be run again to maintain the forecasting error under
519 a fixed value. It can be observed that the error $MSE(M_{n,n+i}(c^{DA}))$ stays
520 smaller than $MSE(M_{n,n+i}(c^n))$ for the 200 time steps shown in Figure 7b,
521 highlighting that the reduction in error gained with the use of DA travels
522 through the model and then benefit positively to the accuracy of the results
523 predicted by Fluidity. In particular, the error $MSE(M_{n,n+31}(c^{DA}))$ at time
524 step $i = 31$ still exhibits a value twice smaller than $MSE(M_{n,n+31}(c^n))$, with
525 a value equal to 1.86.

526 Figure 8 shows a scatter plot of the computed concentrations $M_{n,n+i}(c^n)$
527 and $M_{n,n+i}(c^{DA})$ (using 1391 observations) as a function of the observed
528 data v^{n+i} for i equal to 1, 20, 50, 100 and 200. Figure 8 shows how the

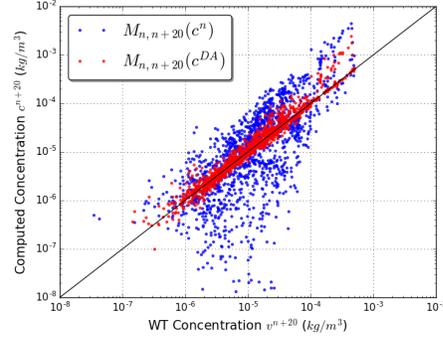
529 pollutant concentration at sensors location evolves through time. It can be
 530 seen that $M_{n,n+i}(c^{DA})$ deviates from the ideal values (black line) starting from
 531 the highest concentrations, i.e. concentrations higher than $1 \times 10^{-4}kg/m^3$,
 532 as shown in Figure 8a. The model tends to recover the higher computed
 533 concentrations very quickly through time. Then, the points having medium
 534 concentration, ranging between $1 \times 10^{-4}kg/m^3$ and $1 \times 10^{-6}kg/m^3$, start to
 535 deviate from the ideal values as i increases, but still keeping a reasonable
 536 spread (Figure 8b, Figure 8c and Figure 8d). Finally, at time step $i = 200$,
 537 the benefit of DA has more or less vanished and the values $M_{n,n+200}(c^{DA})$ tend
 538 to recover the ones obtained from $M_{n,n+200}(c^n)$ (Figure 8e). An interesting
 539 point that can be noted from Figure 8 is that the positive impact and benefit
 540 introduced by the DA process for sensors where the concentration was under-
 541 estimated by Fluidity is preserved through time, i.e the impact is clear even
 542 after 200 time steps (Figure 8e).

543 4.3. Location of assimilating sensors

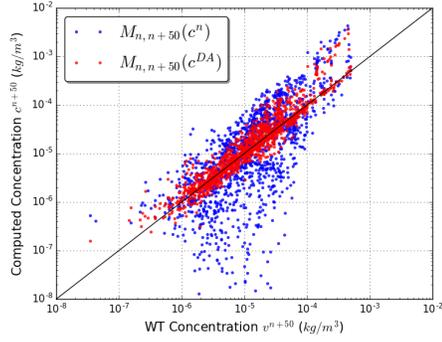
544 The values of the mean squared errors $MSE(c^{DA})$ and $MSE(c^n)$ are here
 545 used to choose an optimal sensors positions which add a positive benefit when
 546 they are assimilated: a trade-off between the number of sensors available in
 547 reality and the gain obtained from the DA process has to be considered.
 548 Several tests were performed to find the optimal sensors positioning and
 549 the cases considered are summarised in Table 2. As a reminder and for
 550 comparison, $MSE(c^n)$ is equal to 3.749.



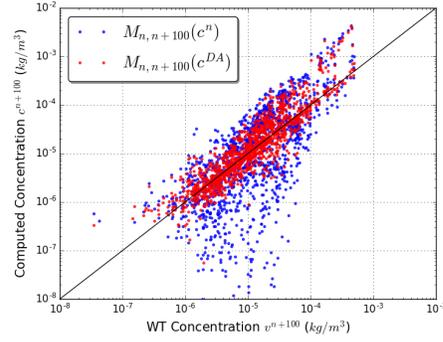
(a) $n + 1$.



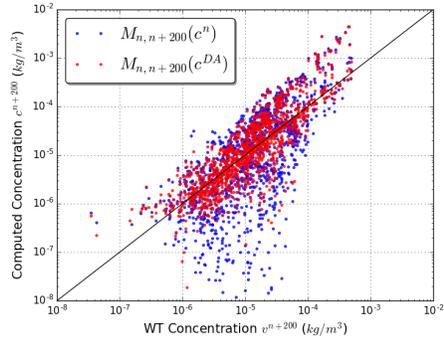
(b) $n + 20$.



(c) $n + 50$.



(d) $n + 100$.



(e) $n + 200$.

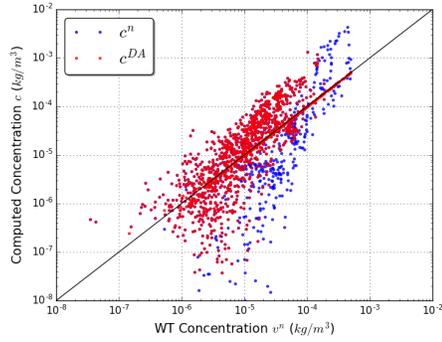
Figure 8: Values of $M_{n,n+i}(c)$ as a function of the wind tunnel data v^{n+i} at time steps (a) $i = 1$, (b) $i = 20$, (c) $i = 50$, (d) $i = 100$ and (e) $i = 200$. The black lines correspond to the ideal matching between data. c^{DA} is obtained assimilating 1391 observations. Logarithm scale is used for all axis.

Test	Sensors location	Number of observations	$MSE(c^{DA})$
Group 1	Figure 4	341	1.13
Group 2	Figure 4	194	3.63
Group 3	Figure 4	300	3.71
Group 4	Figure 4	283	3.73
Group 5	Figure 4	273	3.746
High concentration $c^n \geq 1.5 \times 10^{-4} kg/m^3$	Figure 10a	140	3.89×10^{-1}
Low concentration $c^n \leq 1.5 \times 10^{-6} kg/m^3$	Figure 11a	295	3.746
Group 1 + High Concentration $c^n \geq 1.5 \times 10^{-4} kg/m^3$	Figure 12a	425	3.16×10^{-1}

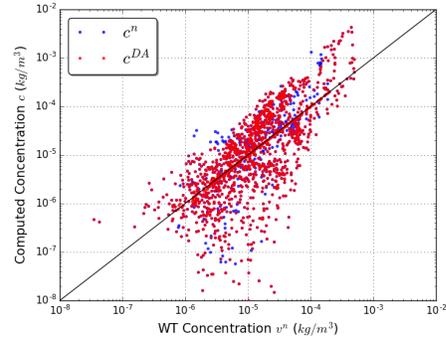
Table 2: Summary of the cases considered to propose an optimal sensor positioning. As a reminder and for comparison, $MSE(c^n)$ is equal to 3.749.

551 As a first attempt, as it can be seen in Figure 4, five different groups of
552 sensors can be identified, based on their distances from the source location.
553 Hence, every sensors group is assimilated separately in order to highlight
554 and quantify their impacts on $MSE(c^{DA})$: results are shown in Table 2 and
555 Figure 9. Assimilating sensors in Group 1 divides the error $MSE(c^{DA})$ by
556 two compared to $MSE(c^n)$, while assimilating any other sensors group lead
557 to a very poor, to not say negligible, reduction in error. As shown in Figure 9,
558 Group 1 is also almost the only one having a positive impact on reducing
559 the spread observed at low concentrations (spread compared to ideal values,
560 i.e. black line). It is then obvious that sensors near the source should be
561 prioritised in order to improve the accuracy of Fluidity.

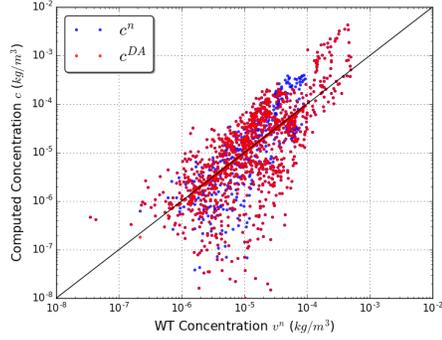
562 The next two tests considered are as follows: only sensors having the high-



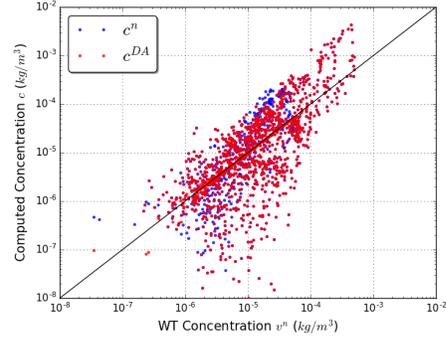
(a) Group 1



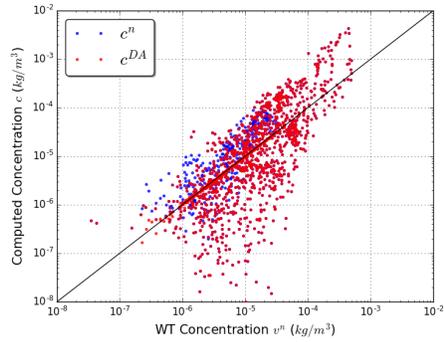
(b) Group 2



(c) Group 3



(d) Group 4

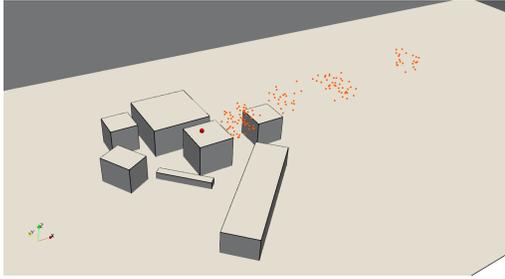


(e) Group 5

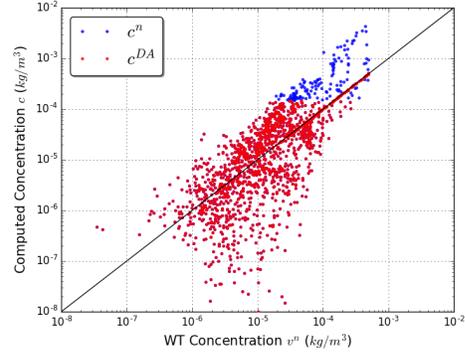
Figure 9: Values of the pollutant concentration c as a function of the wind tunnel data v^n . Value of c^{DA} are obtained assimilating different group of sensors from (a) Group 1, the closest to the source to (e) Group 5, the farthest from the source. For the group labels, see Figure 4. The black lines correspond to the ideal matching between data. Logarithm scale is used for all axis. Data obtained at n , i.e before re-running Fluidity.

563 est Fluidity concentrations ($c^n \geq 1.5 \times 10^{-4} \text{kg/m}^3$) are assimilated, then
 564 only the ones having the lowest concentrations ($c^n \leq 1.5 \times 10^{-6} \text{kg/m}^3$) are
 565 considered. Figure 10a and Figure 11a show the locations of the 140 observa-
 566 tion points having the highest concentrations and the 295 observation points
 567 having the lowest concentrations, respectively. Not surprisingly, the high
 568 concentrations are located downstream and in the alignment of the source.
 569 Results of $MSE(c^{DA})$ are reported in Table 2. The error is reduced by
 570 one order of magnitude if the higher concentrations locations are assimilated
 571 ($MSE(c^{DA}) = 3.89 \times 10^{-1}$), while assimilating the 295 lowest concentrations
 572 locations lead to a non-significant reduction of error ($MSE(c^{DA}) = 3.746$).
 573 Even if the number of assimilated sensors is twice higher, the error is not
 574 significantly reduced when assimilating low concentrations. Looking at the
 575 scatter plots in Figure 10b and Figure 11b, c^n exhibit a large spread around
 576 the ideal value (black line) for low concentrations and assimilating them
 577 sounds a legitimate choice. However, the higher concentrations play a more
 578 determining role in the model error and should then be preferred as location
 579 for sensors. Moreover, talking about air pollution in general, the spots of
 580 high concentration are usually of primary interest and are the ones that need
 581 to be accurately predicted.

582 The last test proposed, as an ultimate optimised case, consists in assimilating
 583 sensors located near the source, i.e. in Group 1, as well as the sensors
 584 exhibiting the highest concentrations only ($c^n \geq 1.5 \times 10^{-4} \text{kg/m}^3$) for all
 585 the other groups. This case leads to 425 observation points as shown in Fig-

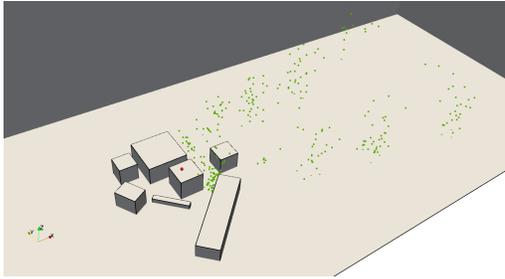


(a)

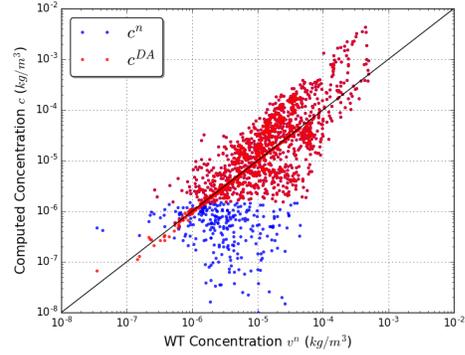


(b)

Figure 10: (a) Location of the sensors assimilated, i.e. the 140 sensors exhibiting the higher Fluidity concentrations ($c^n \geq 1.5 \times 10^{-4} kg/m^3$). The red sphere denotes the location of the source. (b) Values of the pollutant concentration c as a function of the wind tunnel data v^n when the highest concentrations are assimilated. The black line corresponds to the ideal matching between data.



(a)



(b)

Figure 11: (a) Location of the sensors assimilated, i.e. the 295 sensors exhibiting the lower Fluidity concentrations ($c^n \leq 1.5 \times 10^{-6} kg/m^3$). The red sphere denotes the location of the source. (b) Values of the pollutant concentration c as a function of the wind tunnel data v^n when the lowest concentrations are assimilated. The black line corresponds to the ideal matching between data.

586 ure 12b. Table 2 and Figure 12 show the results for this optimised case: the
587 $MSE(c^{DA})$ is equal to 3.16×10^{-1} . Compared to the value obtained when
588 assimilating only the highest concentration location (3.89×10^{-1}), adding
589 sensors in Group 1 in the assimilation process brings a relatively small im-
590 provement of results. However, looking at Figure 12b, this set of sensors
591 positioning remains the optimal one in terms of error reduction: the higher
592 concentrations are properly corrected, thus decreasing the MSE ; while the
593 discrepancies seen at low concentrations are satisfyingly reduced. Hence,
594 the optimal sensors locations recommended to improve the accuracy of Flu-
595 idity results is a trade-off between being close to the source independently
596 of the concentration values and being in region where the concentration is
597 high. In this particular optimal case, 425 observations points are used: as
598 the wind tunnel data are projected onto the *supermesh*, this approximately
599 corresponds to 106 wind tunnel sensors. Compared to the 738 points data
600 provided by the experiment, only 14% of the data need to be used to improve
601 the accuracy of Fluidity.

602 5. Conclusion

603 In this paper an Optimal Three Dimensional Variational (3DVar) data as-
604 simulation model to reduce the discrepancy between CFD results and wind
605 tunnel data in terms of pollutant concentration prediction in urban envi-
606 ronment was presented. Wind tunnel experiments were performed in the

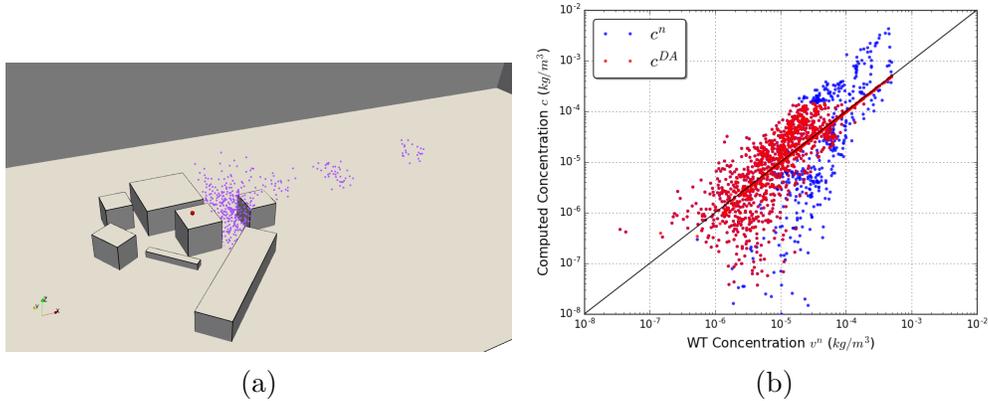


Figure 12: (a) Location of the optimal sensors assimilated, i.e. 425 sensors (Group 1) and sensors having high Fluidity concentrations ($c^n \geq 1.5 \times 10^{-4} kg/m^3$). The red sphere denotes the location of the source. (b) Values of the pollutant concentration c as a function of the wind tunnel data v^n when the optimal sensors are assimilated. The black line corresponds to the ideal matching between data.

607 EnFlo Meteorological Wind Tunnel and the forecasted model was Fluidity,
 608 an open-source CFD software using mesh adaptivity. The mesh adaptivity
 609 technology was used during CFD simulations and then generate an optimal
 610 *supermesh*. The *supermesh* was used in the variational DA process, as well
 611 as to interpolate the Wind Tunnel data.

612 The improvement of Fluidity accuracy, in terms of pollutant concentration
 613 prediction, was discussed using the absolute errors, the mean squared errors
 614 and scatter plots. Using the DA process presented in this paper, the error in
 615 the results between Fluidity and wind tunnel data can be reduced by three
 616 order of magnitudes if all the wind tunnel sensor values are assimilated. It
 617 has been shown that this reduction in error gained using DA is preserved by
 618 the model Fluidity and its benefit can still be observed through several time

619 steps. In particular, it has been observed that high concentration are the one
620 deviating quickly from ideal values, while corrections on low concentrations
621 are fully preserved through time.

622 Finally, an optimal sensors location were proposed taking into account the
623 improvement of Fluidity accuracy while having a limited number of wind
624 tunnel sensors. The optimal sensors locations is a trade-off between being
625 close to the source independently of the concentration values and being in
626 regions where the concentration is high. In the particular case presented in
627 this paper, which used 738 points data from the wind tunnel experiment,
628 only 14% of the data points were needed to reduce the errors by one order
629 of magnitude and improve the accuracy of results predicted by Fluidity in
630 terms of pollutant concentration.

631 **Acknowledgments**

632 This work is supported by the EPSRC Grand Challenge grant Managing Air
633 for Green Inner Cities (MAGIC) EP/N010221/1, by the EPSRC Centre for
634 Mathematics of Precision Healthcare EP/N0145291/1 and by the EPSRC
635 Low Carbon Climate-Responsive Heating and Cooling of Cities (LoHCool)
636 EP/N009797/1.

637 **References**

- 638 [1] World Health Organization. Ambient air pollution: A global as-
639 sessment of exposure and burden of disease. [http://www.who.int/
640 airpollution/ambient/health-impacts/en/](http://www.who.int/airpollution/ambient/health-impacts/en/). Accessed: 2019-04-16.
- 641 [2] United Nations Sustainable Development Goals. Ambient air pollution:
642 A global assessment of exposure and burden of disease. [https://www.
643 un.org/sustainabledevelopment/health/](https://www.un.org/sustainabledevelopment/health/). Accessed: 2019-05-15.
- 644 [3] European Commission. Air quality. [http://ec.europa.eu/
645 environment/air/quality/index.htm](http://ec.europa.eu/environment/air/quality/index.htm). Accessed: 2019-05-15.
- 646 [4] Department of Environment, Food and Rural Affairs, UK. The UK
647 Clean Air Strategy 2019. [https://assets.publishing.service.gov.
648 uk/government/uploads/system/uploads/attachment_data/file/
649 770715/clean-air-strategy-2019.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/770715/clean-air-strategy-2019.pdf). Accessed: 2019-04-16.
- 650 [5] London Environment Strategy. The UK Clean Air Strat-
651 egy 2019. [https://www.london.gov.uk/what-we-do/environment/
652 london-environment-strategy](https://www.london.gov.uk/what-we-do/environment/london-environment-strategy). Accessed: 2019-05-15.
- 653 [6] Y. Tominaga and T. Stathopoulos. CFD simulation of near-field pollu-
654 tant dispersion in the urban environment: A review of current modeling
655 techniques. *Atmospheric Environment*, 79:716–730, 2013.
- 656 [7] M. Lateb, R.N. Meroney, M. Yataghene, H. Fellouah, F. Saleh, and M.C.
657 Boufadel. On the use of numerical modelling for near-field pollutant

- 658 dispersion in urban environments - A review. *Environmental Pollution*,
659 208:271–283, 2016.
- 660 [8] F. Gronwald and S.Y. Chang. Evaluation of the precision and accuracy
661 of multiple air dispersion models. *Journal of Atmospheric Pollution*,
662 6:1–11, 2018.
- 663 [9] N. Mazarakis, E. Kaloudis, A. Nazos, and K.S Nikas. LES and RANS
664 comparison of flow and pollutant dispersion in urban environment. *In-*
665 *ternational Journal of Environmental Studies*, 73:48–65, 2016.
- 666 [10] M.S. Salim, R. Buccolieri, A. Chan, and S. Di Sabatino. Numerical sim-
667 ulation of atmospheric pollutant dispersion in an urban street canyon:
668 Comparison between RANS and LES. *Journal of Wind Engineering and*
669 *Industrial Aerodynamics*, 99:103–113, 2011.
- 670 [11] E. Aristodemou, T. Bentham, C. Pain, R. Colvile, A. Robins, and H. Ap-
671 Simon. A comparison of mesh-adaptive LES with wind tunnel data for
672 flow past buildings: Mean flows and velocity fluctuations. *Atmospheric*
673 *Environment*, 43(39):6238–6253, 2009.
- 674 [12] E. Aristodemou, L.M. Boganegra, L. Mottet, D. Pavlidis, A. Constanti-
675 nou, C. Pain, A. Robins, and H. ApSimon. How tall buildings affect
676 turbulent air flows and dispersion of pollution within a neighbourhood.
677 *Environmental pollution*, 233:782–796, 2018.

- 678 [13] J. Song, S. Fan, W. Lin, L. Mottet, H. Woodward, M. Davies Wykes,
679 R. Arcucci, D. Xiao, J.-E. Debay, H. ApSimon, E. Aristodemou,
680 D. Birch, M. Carpentieri, F. Fang, M. Herzog, G. R. Hunt, R. L. Jones,
681 C. Pain, D. Pavlidis, A. G. Robins, C. A. Short, and P. F. Linden. Nat-
682 ural ventilation in cities: the implications of fluid mechanics. *Building*
683 *Research & Information*, 46(8):809–828, 2018.
- 684 [14] H. Gough, M.-F. King, P. Nathan, C.S.B. Grimmond, A. Robins, C.J.
685 Noakes, Z. Luo, and J.F. Barlow. Influence of neighbouring structures on
686 building façade pressures: Comparison between full-scale, wind-tunnel,
687 CFD and practitioner guidelines. *Journal of Wind Engineering and*
688 *Industrial Aerodynamics*, 189:22–33, 2019.
- 689 [15] L. Mottet, M.-F. King, H. Gough, C.J. Noakes, J.F. Barlow, and C. Pain.
690 Evaluating the capability of a mesh adaptivity LES-CFD software (Flu-
691 idity) to investigate the influence of staggered array on the Silsoe cube
692 façade pressures: a comparative study between two CFD software and
693 Full-Scale experiment. *Journal of Building and Environment.*, —:—,
694 2019.
- 695 [16] E. Kalnay. *Atmospheric Modeling, Data Assimilation and Predictability*.
696 Cambridge University Press, Cambridge, MA, 2003.
- 697 [17] R.E. Kalman. A new approach to linear filtering and prediction prob-
698 lems. *Journal of Basic Engineering.*, 82:35–45, 1960.

- 699 [18] E. Andersson, J. Haseler, P. Undén, P. Courtier, G. Kelly, D. Vasilje-
700 vic, C. Brancovic, C. Cardinali, C. Gaffard, A.Hollingsworth, C. Jakob,
701 P. Janssen, E. Klinker, A. Lanzinger, M. Miller, F. Rabier, A. Simmons,
702 B. Strauss, J-N.Thepaut, and P. Viterbo. The ECMWF implementation
703 of three dimensional variational assimilation (3DVar). Part III: Experi-
704 mental results. *Quarterly Journal of the Royal Meteorological Society.*,
705 124(550):1831–1860, 1998.
- 706 [19] D. M. Baker, W. Huang, Y.R. Guo, J. Bourgeois, and Q.N. Xiao. Three-
707 Dimensional Variational Data Assimilation System for MM5: Implemen-
708 tation and Initial Results. *Monthly Weather Review*, 132:897–914, 2004.
- 709 [20] V. Mons, L. Margheri, J.-C. Chassaing, and P. Sagaut. Data
710 assimilation-based reconstruction of urban pollutant release characteris-
711 tics. *Journal of Wind Engineering and Industrial Aerodynamics*, 169:232
712 – 250, 2017.
- 713 [21] Jorge Sousa, Clara Garca-Snchez, and Catherine Gori. Source appor-
714 tionment and data assimilation in urban air quality modelling for no2:
715 The lyon case study. *Building and Environment*, 132:282 – 290, 2018.
- 716 [22] R. Arcucci, L. Mottet, C. Pain, and Y.-K. Guo. Optimal reduced space
717 for variational data assimilation. *Journal of Computational Physics*,
718 379:51–69, 2019.
- 719 [23] R. Arcucci, C. Pain, and Y.-K. Guo. Effective variational data as-

- 720 simulation in air-pollution prediction. *Big Data Mining and Analytics*,
721 1(4):297–307, 2018.
- 722 [24] Z. Belligoli, R. Dwight, and G. Eitelberg. Assessment of a Data As-
723 simulation Technique for Wind Tunnel Wall Interference Corrections. In
724 *AIAA Scitech 2019 Forum*, page 939, 2019.
- 725 [25] H. Kato and S. Obayashi. Integration of CFD and Wind Tunnel by
726 Data Assimilation. *Journal of Fluid Science and Technology*, 6(5):717–
727 728, 2011.
- 728 [26] D.Q. Zheng, J.K.C. Leung, and B.Y. Lee. An ensemble Kalman filter
729 for atmospheric data assimilation: Application to wind tunnel data.
730 *Atmospheric Environment*, 44:1699–1705, 2010.
- 731 [27] Imperial College London AMCG. Fluidity Manual v4.1.12. https://figshare.com/articles/Fluidity_Manual/1387713, 2015.
732
- 733 [28] P. Gousseau, B. Blocken, T. Stathopoulos, and G.J.F. Van Heijst. CFD
734 simulation of near-field pollutant dispersion on a high-resolution grid:
735 A case study by LES and RANS for a building group in downtown
736 Montreal. *Atmospheric Environment*, 45:428–438, 2011.
- 737 [29] B. Blocken. LES over RANS in building simulation for outdoor and in-
738 door applications: A foregone conclusion? *Building Simulation*, 11:821–
739 870, 2018.

- 740 [30] T. Van Hooff, B. Blocken, and Y. Tominaga. On the accuracy of CFD
741 simulations of cross-ventilation flows for a generic isolated building:
742 Comparison of RANS, LES and experiments. *Building and Environ-*
743 *ment*, 114:148–165, 2017.
- 744 [31] Y. Tominaga and T. Stathopoulos. Ten questions concerning modeling
745 of near-field pollutant dispersion in the built environment. *Building and*
746 *Environment*, 105:390–402, 2016.
- 747 [32] F.T.M. Nieuwstadt, J. Westerweel, and B.J. Boersma. *Turbulence: in-*
748 *roduction to theory and applications of turbulent flows*. Springer, 2016.
- 749 [33] J.W. Deardorff. A numerical study of three-dimensional turbulent
750 channel flow at large Reynolds numbers. *Journal of Fluid Mechanics*,
751 41(2):453–480, 1970.
- 752 [34] J. Bardina, J. Ferziger, and W. Reynolds. Improved subgrid-scale models
753 for large-eddy simulation. In *13th Fluid and Plasma Dynamics Confer-*
754 *ence*, page 1357, 1980.
- 755 [35] T. Bentham. *Microscale modelling of air flow and pollutant dispersion*
756 *in the urban environment*. PhD thesis, Imperial College London, 2004.
- 757 [36] C. Pain, A.P. Umpleby, C.R.E. De Oliveira, and A.J.H. Goddard. Tetra-
758 hedral mesh optimisation and adaptivity for steady-state and transient
759 finite element calculations. *Computer Methods in Applied Mechanics*
760 *and Engineering*, 190:3771–3796, 2001.

- 761 [37] D. Pavlidis, G.J. Gorman, J.L.M.A. Gomes, C. Pain, and H. ApSi-
762 mon. Synthetic-Eddy Method for Urban Atmospheric Flow Modelling.
763 *Boundary-Layer Meteorology*, 136:285–299, 2010.
- 764 [38] N. Jarrin, S. Benhamadouche, D. Laurence, and R. Prosser. A synthetic-
765 eddy-method for generating inflow conditions for large-eddy simulations.
766 *International Journal of Heat and Fluid Flow*, 27(4), 2006.
- 767 [39] S.B. Pope. *Turbulent Flows*. Cambridge University Press, 2000.
- 768 [40] H. K. Engl, M. Hanke, and A. Neubauer. Regularization of Inverse
769 Problems. *Kluwer*, 1996.
- 770 [41] N. Nichols. *Data Assimilation - Chapter Mathematical concepts in data*
771 *assimilation*. Springer, 2010.
- 772 [42] P.C. Hansen. *Rank Deficient and Discrete Ill-Posed Problems*. Society
773 for Industrial and Applied Mathematics, 1998.
- 774 [43] L. D’Amore, R. Arcucci, L. Marcellino, and A. Murli. A Parallel Three-
775 dimensional Variational Data Assimilation Scheme. In *AIP Conference*
776 *Proceedings*, volume 1389, pages 1829–1831, 2011.
- 777 [44] A.C. Lorenc. Development of an operational variational assimilation
778 scheme. *Journal of the Meteorological Society of Japan*, 75:339–346,
779 1997.

- 780 [45] J.P. Courtier. A strategy for operational implementation of 4D-VAR,
781 using an incremental approach. *Quarterly Journal of the Royal Meteorological Society*, 120(519):1367–1387, 1994.
- 783 [46] E.N. Lorenz. Empirical orthogonal functions and statistical weather
784 prediction., 1956.
- 785 [47] A. Hannachi, I.T. Jolliffe, and D.B. Stephenson. Empirical orthogonal
786 functions and related techniques in atmospheric science: A review. *International Journal of Climatology: A Journal of the Royal Meteorological Society*, 27:1119–1152, 2007.
- 789 [48] A. Hannachi. A Primer for EOF Analysis of Climate Data. *Department of Meteorology, University of Reading, UK*, 2004.
- 791 [49] R. Arcucci, L. D’Amore, J. Pistoia, R. Toumi, and A. Murli. On the
792 variational data assimilation problem solving and sensitivity analysis. *Journal of Computational Physics*, pages 311–326, 2017.
- 794 [50] AG Robins. Experimental model techniques for the investigation of the
795 dispersion of chimney plumes. *Proceedings of the Institution of Mechanical Engineers*, 189(1):44–54, 1975.
- 797 [51] W.J. Coirier and K. Sura. CFD modeling for urban area contaminant
798 transport and dispersion: model description and data requirements. In *Sixth Symposium on the Urban Environment, The 86th AMS annual meeting*, volume 163, pages 175–185, 2006.
- 800

- 801 [52] J. Franke, A. Hellsten, H. Schlnzen, and B. Carissimo. *Best practice*
802 *guideline for the CFD simulation of flows in the urban environment,*
803 *COST Action 732.* COST Office, 2007.
- 804 [53] M. Carpentieri, P. Salizzoni, A. Robins, and L. Soulhac. Evaluation of
805 a neighbourhood scale, street network dispersion model through com-
806 parison with wind tunnel data. *Environmental Modelling & Software,*
807 *37:110–124,* 2012.
- 808 [54] V. Fuka, Z.-T. Xie, I.P. Castro, P. Hayden, M. Carpentieri, and
809 A. Robins. Scalar Fluxes Near a Tall Building in an Aligned Array
810 of Rectangular Buildings. *Boundary-Layer Meteorology,* 167(1):53–76,
811 2018.