

Stochastic perturbations and dimension reduction for modelling uncertainty of atmospheric dispersion simulations

Sylvain Girard^{a,*}, Patrick Armand^b, Christophe Duchenne^b, Thierry Yalamas^c

^a18 boulevard de Reuilly, Phimeca, France

^bCEA, DAM, DIF, F-91297 Arpajon, France

^cCentre d'affaires du Zénith, 34 rue de Sarliève, 63800, Cournon d'Auvergne, France

Abstract

1 Decision of emergency response to releases of hazardous material in the at-
2 mosphere increasingly rely on numerical simulations. This paper presents
3 two contributions for accounting for the uncertainty inherent to those simu-
4 lations. We first focused on one way of modelling these uncertainties, namely
5 by applying stochastic perturbations to the inputs of the numerical dispersion
6 model. We devised a generic mathematical formulation for time dependent
7 perturbation of both amplitude and dynamics of the inputs. It allows a more
8 thorough exploration of possible outcomes than simpler perturbations found
9 in the literature. We then improved on the current state of the art on di-
10 mension reduction of atmospheric data. Indeed, most statistical methods
11 cannot cope with high dimensional data such as the maps simulated with at-
12 mospheric dispersion models. Principal component analysis, the most widely
13 used method for dimension reduction, relies on a linearity hypothesis that is
14 not verified by these sets of maps. We conducted a very encouraging exper-
15 iment with auto-associative models, a non-linear extension of this method.

Keywords: atmospheric dispersion, uncertainty propagation, wind field, time warp, perturbation

*Corresponding author

Email addresses: `girard@phimeca.com` (Sylvain Girard), `patrick.armand@cea.fr` (Patrick Armand), `christophe.duchenne@cea.fr` (Christophe Duchenne), `yalamas@phimeca.com` (Thierry Yalamas)

16 1. Introduction

17 Decision of emergency response to releases of hazardous material in the
18 atmosphere increasingly rely on predictions from numerical models. Such
19 simulations of atmospheric dispersion are highly uncertain due to the com-
20 plexity of the physical phenomena, and because their inputs, in particular
21 meteorological or source term related, are highly uncertain. We propose two
22 methodological improvements to the current practices aimed at accounting
23 for these uncertainties.

24 In section 2, we state the importance of setting the decision problem in
25 a probabilistic framework, and introduce a realistic case study later used to
26 illustrate our two contributions.

27 In section 3 we expound on one way of modelling these uncertainties,
28 namely by applying stochastic perturbations to the inputs of the numerical
29 dispersion model. We devised a generic mathematical formulation for time
30 dependent perturbation of both amplitude and dynamics of the inputs. It
31 allows a more thorough exploration of possible outcomes than simpler per-
32 turbations found in the literature.

33 The output of dispersion models are spatial maps. Analysing a set of
34 maps, whether qualitatively by visual inspection, or quantitatively with sta-
35 tistical methods, is much harder than dealing with numerical values. In
36 section 4 we discuss the issue of obtaining a concise representation of maps
37 by a few scalars. Principal component analysis is the most widely used
38 method for dimension reduction. It relies however on a linearity hypothe-
39 sis that is seldom verified by sets of maps produced by dispersion models.
40 We conducted an encouraging experiment with auto-associative models, a
41 non-linear extension of this method.

42 Both sections 3 and 4 begin with a detailed survey of the literature on
43 the topic at hand.

44 2. Problem statement

45 We consider the following idealised decision problem. Hazardous material
46 is released in the atmosphere during a given period of time. Mitigation ac-
47 tions, for instance population sheltering or evacuation, must be performed in
48 areas where a given concentration threshold is exceeded. These concentra-
49 tions are predicted with a physical model simulating transport and dispersion
50 in the atmosphere, and deposition by rain.

This decision making scenario is inspired by a real industrial incident that happened on January 2013 at the Lubrizol chemical plant located in Rouen, France. Operation mistakes and minor system failures in the plant resulted in extended release of hydrogen sulphur and mercaptan, which are both foul-smelling. The first report of olfaction in the neighbourhood of the site occurred at 8:00 am (local time) on Monday 21 January 2013. The first major emission peak occurred 12 hours later (2013-01-21 20:00). The major part of the material inventory (99%) was emitted during 23 hours, between 2013-01-21 13:30 and 2013-01-22 12:30 [19]. The wind blew the plume as far as Paris during Monday night and towards London on Tuesday. Thousands of people have complained of nausea and headaches. For practical reasons, we focused here on a restricted area of about 35 kilometres horizontal print. Our approach could equally be applied at different space scales.

2.1. *Physical model*

The dispersion simulations were carried out with Parallel-Micro-SWIFT-SPRAY (PMSS). Originally, Micro-SWIFT-SPRAY (MSS) [33] was developed in order to provide a simplified but rigorous computational fluid dynamics solution of the flow and dispersion over rugged terrains and built-up environments in a limited amount of time. MSS encompasses the local scale high resolution versions of the SWIFT and SPRAY models.

SWIFT is a 3D terrain-following mass-consistent diagnostic model taking account of the buildings and providing the 3D fields of wind, turbulence, and temperature. SWIFT interpolates between meteorological measurements (ground stations and vertical profiles), numerical data issued by meso-scale simulations (as in this paper) and, possibly, analytical relations of the flow influenced by the buildings (displacement zone, wake zone, skimming zone, etc.).

SPRAY is a 3D Lagrangian Particle Dispersion Model able to account for the presence of buildings. Both SWIFT and SPRAY can deal with complex terrains and evolving meteorological conditions and with specific features of the release (heavy gas, light gas, etc.). More recently, SWIFT and SPRAY have been efficiently parallelized in time, space, and numerical particles leading to the PMSS modelling system [28]. PMSS has been thoroughly validated against several wind tunnel and in-field experimental campaigns in the framework of notably the European COST Action [34] and the UDINEE project [27]. The performances of PMSS give full satisfaction as the modelling system

87 is compliant with the validation criteria for 3D dispersion models adapted to
 88 built-up areas, proposed by Hanna and Chang [16] and used internationally,

89 2.2. Deterministic decision map

90 The source is located in the middle of the simulation domain, paral-
 91 lelepiped whose horizontal print is a square with edge of 35 km. The simula-
 92 tion duration was set 35 hours in order to ensure that all material has either
 93 been deposited or exited the simulation domain at the end of simulation. In
 94 a deterministic framework, a single simulation is run using the most credi-
 95 ble values for the meteorological and source term model inputs. From now
 96 on, we call this set of values the (input) conjecture, and likewise we refer to
 97 the concentrations simulated using them as conjectured concentrations. We
 98 focus here on three uncertain inputs known to have a substantial impact on
 99 simulation output [3, 14, 2, 15]:

- 100 • the rate of emission of material is a time series, called source term,
- 101 • the rain intensity is a scalar spatio-temporal field,
- 102 • and wind velocity (speed and direction) is a vector spatio-temporal
 103 field.

104 The conjectured source term, displayed on figure 1, was adapted from
 105 data established by Ismert and Durif [20].

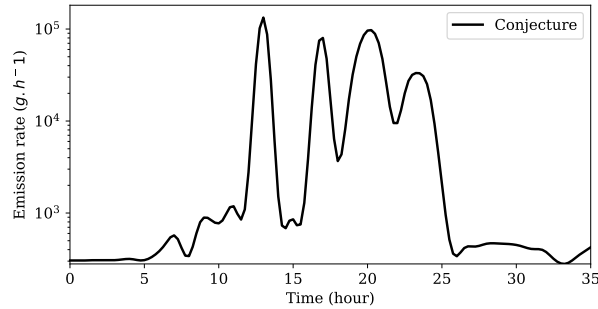


Figure 1: Conjectured source term. The time abscissa starts on 21 January 2013, 7:00.

106 The conjectured rain and wind fields were obtained from the community
 107 reconstruction weather and forecast meso-scale modelling system WRF [31].
 108 The WRF simulation domain has a horizontal resolution of 1 km. For the

109 wind, we used a set of 514 vertical profiles of the horizontal components, and
 110 we kept the 21 vertical layers below 3 km above ground level (AGL), plus
 111 surface data at 10 m AGL. Their locations are displayed in figure 2. WRF
 112 simulations are sampled every 15 minutes, and we used the 141 time steps
 113 from 21/01/2013 07:00 to 22/01/2013 18:00.

114 We want to predict, at each location, whether an arbitrary concentra-
 115 tion threshold of $2\mu g \cdot cm^{-3}$ is exceeded during the considered 35 hour time
 116 frame. In real application, this decision criterion could correspond for in-
 117 stance to olfaction threshold or important health damage. The area where
 118 the concentration threshold is exceeded during the conjecture simulation is
 119 coloured in dark red in figure 2. Note that in this example, we do not con-
 120 sider the actual olfactory limit of the mercaptan which is lower than the
 121 chosen threshold and would lead to a decision map spanning quite all over
 122 the simulation domain, thus making less interesting the visual presentation
 123 of the uncertainties influence on the decision map.

124 2.3. Probabilistic decision map

125 We tackle the decision problem stated above from the viewpoint of un-
 126 certainty propagation. We need the probability distribution of maximum
 127 concentrations, but cannot model it directly. Instead, we model the un-
 128 certainty of inputs of the physical model by stochastic perturbations. The
 129 atmospheric dispersion model is a deterministic function. Here, its inputs are
 130 the source term, and the rain and wind fields. Its output is the maximum
 131 concentration over time at each location. Figure 3 represents the chain of
 132 functions linking the random variables of the problem. The random vector
 133 of perturbed inputs is denoted by Y , and the simulated maximum concen-
 134 tration over time by $Z(s)$, a function of the location s . The upstream part
 135 of the chain will be described in section 3.

136 In a probabilistic framework, taking a decision implies to admit a risk
 137 of committing a specified error. We choose here the risk of deciding not
 138 to perform a mitigation action while the concentration threshold is actually
 139 exceeded. The corresponding probabilistic decision rule is to decide action
 140 where the estimated probability of exceedance is above a small arbitrarily
 141 specified value, for instance 5%. The other possible error, not considered
 142 here, would be to decide unnecessary actions.

143 Exceedance maps are obtained from maximum concentration maps by
 144 setting each cell to 1 if the chosen threshold is exceeded and 0 otherwise.
 145 Averaging a set of exceedance maps yields a Monte Carlo estimator of the

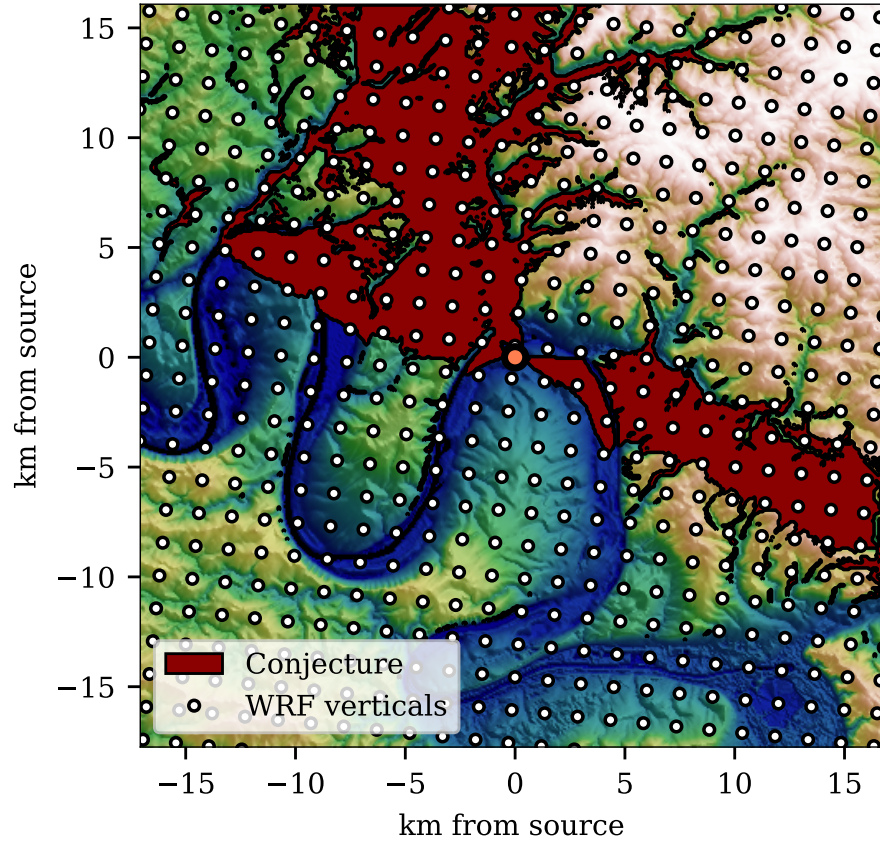


Figure 2: Area where the concentration threshold is exceeded during the conjecture simulation (dark red). The source location is indicated by an orange dot. White dots indicate the locations of the WRF verticals used as meteorological conjecture.

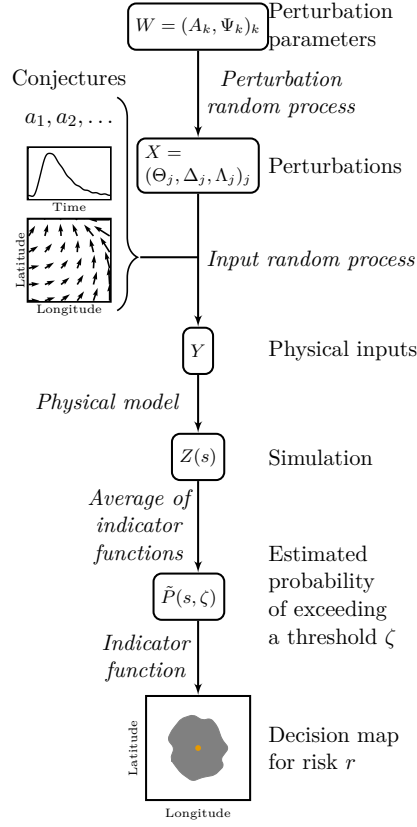


Figure 3: Function chain of the uncertainty propagation. Capital framed letters designate random vectors. Arrows annotated in italic designate functions.

146 probability of exceedance. Figure 4 displays probability of exceedance thus
 147 estimated from a sample of 100 simulations with inputs perturbed as de-
 148 scribed in section 3. This relatively small sample size is representative of
 149 crisis context when decision must be taken rapidly while each simulation
 150 requires up one hour of computation as in our case study. Assessing the
 151 convergence of the estimator of small probability with small samples is not
 152 trivial [1], but we expect errors due to partial convergence. Therefore, a
 153 conservative stance is to draw a decision boundary enclosing the level set
 154 corresponding to the chosen probability threshold (for instance the black line separating green and yellow in figure 4, for a 5% threshold).

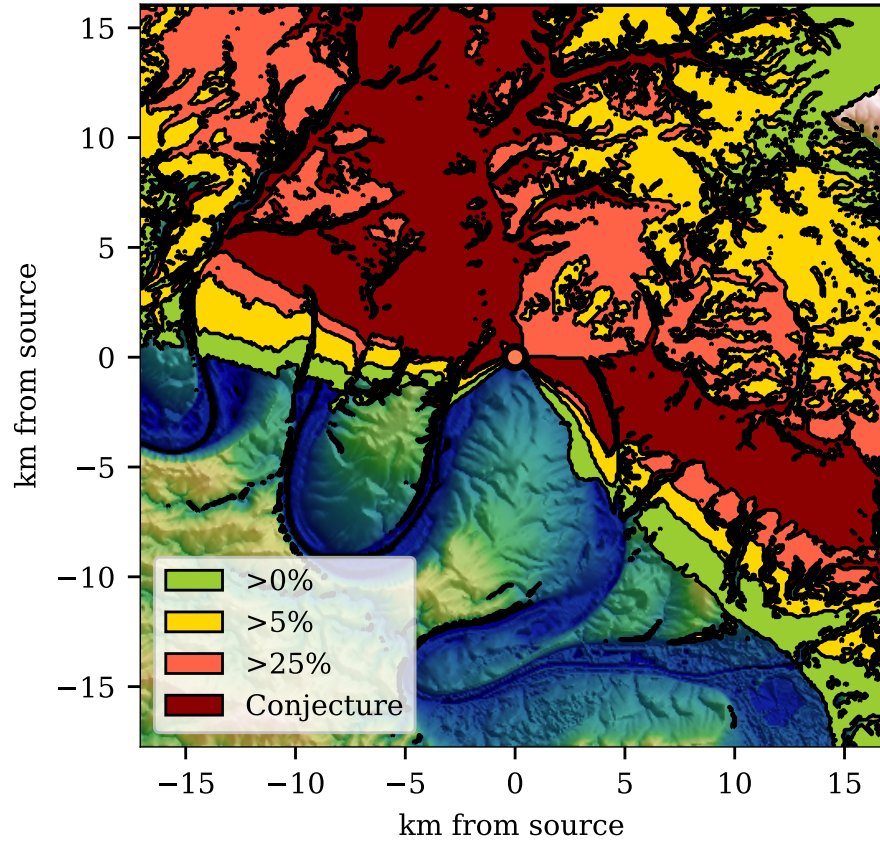


Figure 4: Orange, yellow and green areas correspond to three intervals for the probability of exceedance estimated from 100 simulations.

156 3. Amplitude and dynamics perturbation scheme

157 We will now describe specifically the perturbation model. In section 3.1,
158 we review the rationale behind modelling uncertainty by stochastic pertur-
159 bations and the state of the art. In section 3.2, we present our generic math-
160 ematical formulation for amplitude and dynamics perturbations. Finally, in
161 section 3.3 we detail how it was implemented for the case study.

162 3.1. Uncertainty modelling by stochastic perturbations

163 The input conjecture comprises any model input that is issued from par-
164 tial or imprecise observations, or from other physical model simulations, and
165 possibly also model parameters that implicitly account for the unavoidable
166 discrepancy between the real complex system under study and its idealised
167 mathematical representation.

168 One possible approach to modelling those uncertainties is to apply stochas-
169 tic perturbations to the input conjecture. The meteorological conjecture is
170 usually obtained from simulations with meteorological models. Several au-
171 thors suggested to apply uncertainty propagation to those models, or to use
172 sets of different models to produce build ensembles of input meteorological
173 conditions [35, 10]. However, as [7] points out, it is rather unlikely that such
174 an ensemble of conjectures might be available in a crisis context. Indeed,
175 considerable efforts must be spent to calibrate such ensembles [11, 36], and
176 the calibration process requires reference data that are not available when
177 dealing with accidental releases. It can be expected that using ensembles
178 designed for meteorological forecast as a substitute for specifically calibrated
179 ones would result in underestimation of uncertainty. However ensemble ap-
180 proach and conjecture perturbations are not mutually exclusive: applying
181 additional perturbations to the ensemble members could possibly retrieve
182 the missing variability [23]. We focus here on the case when a single con-
183 jecture is available.

184 Random perturbations of inputs commonly found in the literature [7, 4,
185 17, 12, 14, 15] are time independent random variables. They often follow
186 a Gaussian distribution for additive perturbation, and log-normal for multi-
187 plicative perturbation. The dynamics of the conjecture is rarely perturbed,
188 except sometimes by global time delays [14, 15].

189 3.2. General mathematical formulation

190 The conjecture is, in general, a set of data of diverse dimensions: scalars,
191 time series, and spatio-temporal fields. They are grouped by a brace on

192 figure 3. The random perturbations, collectively denoted by X on figure 3,
 193 are functions of the set of random perturbation parameters W .

194 More precisely, let $Y(s, t)$ be the spatio-temporal random vector obtained
 195 by perturbation of a conjecture $c(s, t)$. We adopted the following generic
 196 perturbation:

$$Y(s, t) = \Gamma(s, \Theta(t)) c(s, \Theta(t)) + \Delta(s, \Theta(t)). \quad (1)$$

197 The random vector Γ (respectively Δ) is a multiplicative (respectively addi-
 198 tive) perturbation of the amplitude of the conjecture, that we will call the
 199 “gain” (respectively “offset”) perturbation. The random function Θ is the
 200 perturbation of the dynamics of the conjecture, called “time warp” pertur-
 201 bation.

202 3.2.1. Spatio-temporal structure of the perturbation

203 The spatio-temporal structure of each component of the perturbation
 204 must be postulated. We chose to impose smooth oscillating temporal vari-
 205 ations, and have the perturbation depend on the location only through the
 206 conjecture. More precisely, we used sums of cosines

$$\sum_{k=1}^K A_k \cos(2\pi\omega_k t + \Psi_k). \quad (2)$$

207 with random phases Ψ_1, \dots, Ψ_K , and random amplitudes A_1, \dots, A_K . The
 208 time structure of the random process is controlled by the choice of the K
 209 periods $\omega_1, \dots, \omega_K$. The value of K is itself a parameter to be chosen by the
 210 modeller.

211 The phases are independent and uniformly distributed on $[0, 2\pi]$. The
 212 distributions of the amplitudes are arbitrary.

213 The gain and offset random processes are directly given by equation (2).
 214 The additional derivation of the time warp process is detailed in the next
 215 section.

216 3.2.2. Time warp

217 Let $\phi : \mathbb{R}^+ \times \mathbb{R}^+ \mapsto \mathbb{R}^+$ be a function that expands or contracts a time
 218 interval δt by a factor $\lambda(t)$ that varies in time:

$$\phi : t, \delta t \mapsto \phi(t, \delta t) = \lambda(t) \delta t. \quad (3)$$

219 Let $\{0 = t_0 < t_1 < \dots < t_q = T\}$ be a sequence of instants interspersing a
 220 time frame of duration T . Denote by $\phi_0, \dots, \phi_{q-1}$ the warped time intervals:

$$221 \quad \phi_n = \phi(t_n, t_{n+1} - t_n). \quad (4)$$

222 The associated *time warp* function $\theta : \mathbb{R}^+ \mapsto \mathbb{R}^+$ preserves the time origin
 223 and warps subsequent instants:

$$\theta(t_0) = t_0 = 0, \quad (5)$$

224 and $\forall n : 1 \leq n \leq q - 1$,

$$\theta(t_n) = \frac{T}{\sum_{i=0}^{q-1} \phi_i} \sum_{i=0}^n \phi_i. \quad (6)$$

225 It follows from the above definition that a time warp function also preserves
 226 the total duration: $\theta(t_q) = t_q = T$.

227 A time warp random process Θ is fully characterised by specifying a ran-
 228 dom process Λ for generating warping factors $\lambda(t)$. We used smoothly oscil-
 229 lating functions as defined by equation (2). In practice, a warped time series
 230 (namely a realisation $y(s, t)$ of the random function $Y(s, t)$ in equation (1))
 231 is obtained by

- 232 1. Applying the gain and offset perturbations to the conjecture.
- 233 2. Sampling warping factors $\lambda(t)$ from Λ .
- 234 3. Computing the corresponding warped instants $\theta(t_0), \theta(t_1), \dots, \theta(t_q)$.
- 235 4. Interpolating the time series resulting from step 1 at the warped in-
 236 stants.

237 3.3. Perturbations for the case study

238 We perturbed four input conjectures by processes based on equation (1):

- 239 • source term (gain and time warp),
- 240 • rain intensity (gain and time warp),
- 241 • wind speed (offset and time warp),
- 242 • and wind direction (offset and time warp).

Component	Criterion
All time warps	95% of the values lag behind (or anticipate) the corresponding conjectured value by less than 2 hours.
Source term gain	95% of the gain factors (as time varies, and from one realisation to another) are between 0.5 and 2, with a median equal to 1.
Rain intensity gain	95% of gain factors (as time varies, and from one realisation to another) are between 0.5 and 2, with a median equal to 1.
Wind speed and direction offsets	95% of offsets lie within an interval whose length depends on the spatial average of wind speed in each vertical layer at each time step.

Table 1: Calibration criteria for the amplitude distributions.

243 This amounts to a total of 8 random functions (2 gains, 2 offsets and 4 time
244 warp). All these components of the perturbation processes are statistically
245 independent. They have all the same temporal structure, with $K = 3$ periods.
246 Each component thus involves 6 random variables: 3 phases and 3 amplitudes
247 (see equation (2)). The total number of random variables is $8 \times 6 = 48$.

248 We choose the following values for the periods: $\omega_0 = 0$, $\omega_1 = T/4$ and
249 $\omega_3 = T$ (the simulation time frame T is equal to 35 h). They induce contri-
250 butions to the perturbation that are respectively, constant in time, with 4
251 cycles within the time frame, and with a single cycle.

252 We used independent Gaussian random variables with zero mean for all
253 amplitudes (denoted A_k above). We applied an exponential transformation to
254 the gain perturbations. The resulting distributions of multiplicative factors
255 are thus roughly log-normal with median equal to 1. The standard deviations
256 were determined following the criteria listed in table 1 established by expert
257 judgement. Refer to supplementary material for illustrations.

258 4. Dimension reduction of a set of maps

259 Precise estimation of small probabilities of exceedance require larger sam-
260 ple size than what can usually be achieved in a crisis context, due to the sub-
261 stantial CPU cost of detailed atmospheric dispersion models. Furthermore,
262 input uncertainty models rely on many postulates, for instance the distri-
263 bution of perturbation parameters. The robustness of the decision criterion

can be tested by repeating the uncertainty propagation with different sets of perturbation parameters, which would require even greater sample size. Sample size can be a limiting factor even in less time constrained contexts like sensitivity analysis [14, 15] and source term estimation by inverse methods (Liu, 2017).

Model emulation is an alternative to direct estimation by Monte Carlo sampling [3, 2, 15, 26, 24, 25]. The simulation sample is used to build a mathematical approximation of the physical model whose computational cost is negligible. Emulation techniques, such as Gaussian process regression [30] apply to models with a scalar output, but the output of the physical model considered here is a spatial map. In practice it is represented by a large number of values sampled on the nodes of a grid, in the same manner that a raster image is a set of pixels. Because these node variables are intricately dependent on one another, we can attempt parametrising the maps with a lesser number of variables by a process similar to those used for image compression.

Principal component analysis (PCA) is a method for dimension reduction used extensively in data analysis for more than a century [29, 21]. It can be considered as the state of the art for dimension reduction in the specific field of dispersion simulation [6, 32, 26, 24, 25]. However, it relies on a linearity hypothesis that is not verified by the set of maps typically produced by atmospheric dispersion models. We expound on this issue in section 4.1.

Auto-associative models (AAM) are a non-linear extension of PCA that benefits from an explicit and attractive mathematical foundation [13]. It is presented in section 4.2, and its performance is compared to that of PCA in section 4.3.

4.1. Limitations of principal component analysis

The output of the dispersion model is a spatial map (denoted by $Z(s)$ in figure 3) discretised on a grid with p nodes. As such, it can be seen as a p dimensional vector, and the set of output maps \mathcal{Z} is a subset of \mathbb{R}^p . The principle of dimension reduction is to build an approximate low dimensional system of coordinates for \mathcal{Z} based on a sample of elements z_1, \dots, z_N . We will assume without loss of generality that the sample point cloud is centred on the origin. Otherwise, the procedures described in the following must simply be preceded by a translation of vector $\pm \frac{1}{N} \sum_{i=1}^N z_i$

The algorithm of PCA solves a sequence of optimisation problems:

300 For $k = 1, \dots, p$, find the unit vector a_k orthogonal to any a_i with $i < k$
 301 that maximises $\sum_{i=1}^n (a'_k z_i)^2$ (where a'_k denotes the transpose of a_k).

302 For any dimension d , the principal directions a_1, \dots, a_d form an orthonor-
 303 mal basis of a linear space \mathcal{L}_d approximating \mathcal{Z} . The low dimensional coordi-
 304 nates are simply the coordinates in this basis. We call *residual* the difference
 305 $z - \sum_{i=1}^d (a'_i z) a_i$ between an element of \mathcal{Z} and its approximation.

306 It can be shown that \mathcal{L}_d is such that

- 307 • the sum of squares of the sample residuals is minimised,
- 308 • the sum of squares of the Euclidean distances between sample points
 309 is best preserved by projection.

310 A major advantage of PCA is that those equivalent optimisation problems
 311 have a closed form solution: the (a_i) are the eigen vectors of $\frac{1}{N} Z'Z$, where Z
 312 is the $N \times p$ matrix whose rows are the (z_i) .

313 PCA is efficient, in the sense that it yields good low dimensional ap-
 314 proximations, when the relations between the p original variables are linear.
 315 Unfortunately, this is seldom verified when the p original variables are sam-
 316 pled on a time series or a spatial map.

317 Consider for instance a set of time series consisting of a single identical
 318 bell shaped pulse occurring at varying instants. Discretising it at p evenly
 319 spaced abscissas yields a p dimensional point cloud as before. This sim-
 320 ple example devised by Fukunaga and Olsen [9] is actually representative of
 321 many situation common in atmospheric dispersion: time evolution of a con-
 322 centration when a plume passes over a recording station, or comparison of
 323 circular cross sections of plumes with differing orientations. The minimum
 324 number of parameters needed to reversibly encode a data set without loss of
 325 information is often referred to as its “intrinsic dimension” [18]. Here, it is
 326 equal to 1. Indeed, a single scalar, say the abscissa of the top of the bell,
 327 is sufficient to fully parametrize the set of curves. However, the number of
 328 principal directions required to achieve a good approximation is close to p .
 329 In a previous communication [8], we applied PCA to sets of dispersion sim-
 330 ulations resulting from low dimensional perturbations. While the first two
 331 or three principal components carried much more information than the sub-
 332 sequent ones, at least a dozen of them were required to properly encode the
 333 original dataset. This lack of efficiency even with simple perturbations shows
 334 that PCA is not well suited for analysing the output of complex perturbation
 335 schemes.

336 4.2. Non linear dimension reduction with auto-associative models

337 The algorithm for building an AAM starts with the $N \times p$ sample data
 338 matrix $Z_1 = Z$ whose rows are the (z_i) , and repeats the following steps for
 339 $k = 1, \dots, p$:

- 340 1. Find a direction a_k minimising a loss function.
- 341 2. Compute the vector of projection coordinates $c_k = Z_k$.
- 342 3. Estimate the recovery function r_k , namely an approximation of the
 343 function linking the components of c_k to the rows of Z_k .
- 344 4. Set $Z_{k+1} = Z_k - \tilde{Z}_k$, where \tilde{Z}_k is the $N \times p$ sample data matrix whose
 345 rows are the images of the projection coordinates by r_k .

346 Following the recommendations of Girard and Iovleff [13], we used a loss
 347 function that best preserves nearest neighbours, and built the recovery (r_k)
 348 with cubic splines. Note that PCA is a linear special case of AAM with the
 349 loss function given in the previous section, and the linear maps $r_k : \alpha \mapsto \alpha a_k$
 350 for recoveries.

351 It was shown theoretically that AAM surpasses neural networks with
 352 simple architectures such as auto-associative perceptron with one hidden
 353 layer [13]. More sophisticated networks seem capable of good performances
 354 [22], but they require large training samples. AAM is less inductive, but its
 355 added rigidity makes it able to cope with small training samples. In that
 356 respect, our approach is closer to that of Bowman and Woods [5].

357 4.3. Compared performances of PCA and AAM

358 An important feature that one expects from a good dimension reduction
 359 technique is the fidelity of the projected data to the original. The quadratic
 360 mean of the residuals (the difference between projected and original data) is
 361 a common measure of the missing information in the projection. As a matter
 362 of fact, principal components are the solution of the minimisation of this very
 363 quantity. It is also common to normalise this measure by the variance of the
 364 original data. Indeed, when the data are realisations of random variables, the
 365 quadratic mean of the residuals is an estimate of the associated variance. This
 366 allows computing the amount of the overall variance that is explained by each
 367 projection direction. Following Girard and Iovleff [13], we call “information
 368 ratio” the sum of the variances explained by a set of directions.

Figure 5 compares the information ratios of the AAM and PCA projections of increasing dimension. It shows that the first few AAM directions are much more informative than that of PCA. Indeed, AAM is able to account for almost 80% of the data set variance with two parameters only, while PCA needs six directions to catch up with this value.

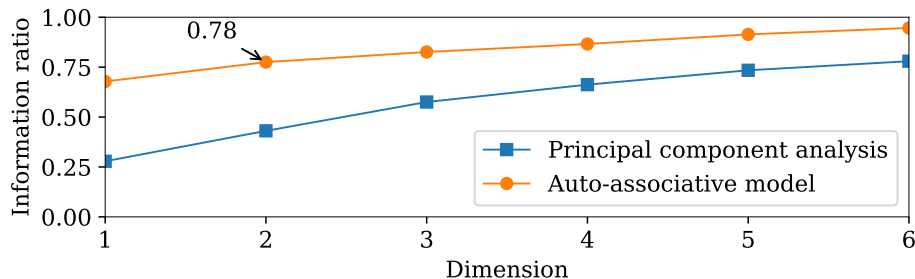


Figure 5: Information ratio of the AAM and PCA projections of increasing dimension.

In the following paragraphs, we compare in more details the AAM and PCA 2 dimensional projections. We chose to leave aside subsequent directions for three reasons:

- 2D projections can be plotted and are therefore better suited for our illustrative purpose.
- The limited size of the training sample, 100 simulations, induces a significant risk of overfitting. Cross-validation showed indeed that the third AAM direction is unreliable for extrapolation.
- Very low dimensional projection (3D and below) are the most versatile. Many methods stop working in dimension 4 and above.

Each row of figure 6 shows a group of three simulated maps whose projections are close to one another in the 2D AAM coordinate system. They are labelled by coloured letters that locate them the scatter plot of the 2D AAM coordinate system in figure 7. The simulated maps within a group are similar, while maps from different groups are dissimilar. Each groups can be characterised by the main features of the exceedance area (contoured in pale red), namely an extension towards the south-east direction, and the direction and width of the northern fan shaped area. These observations suggest that

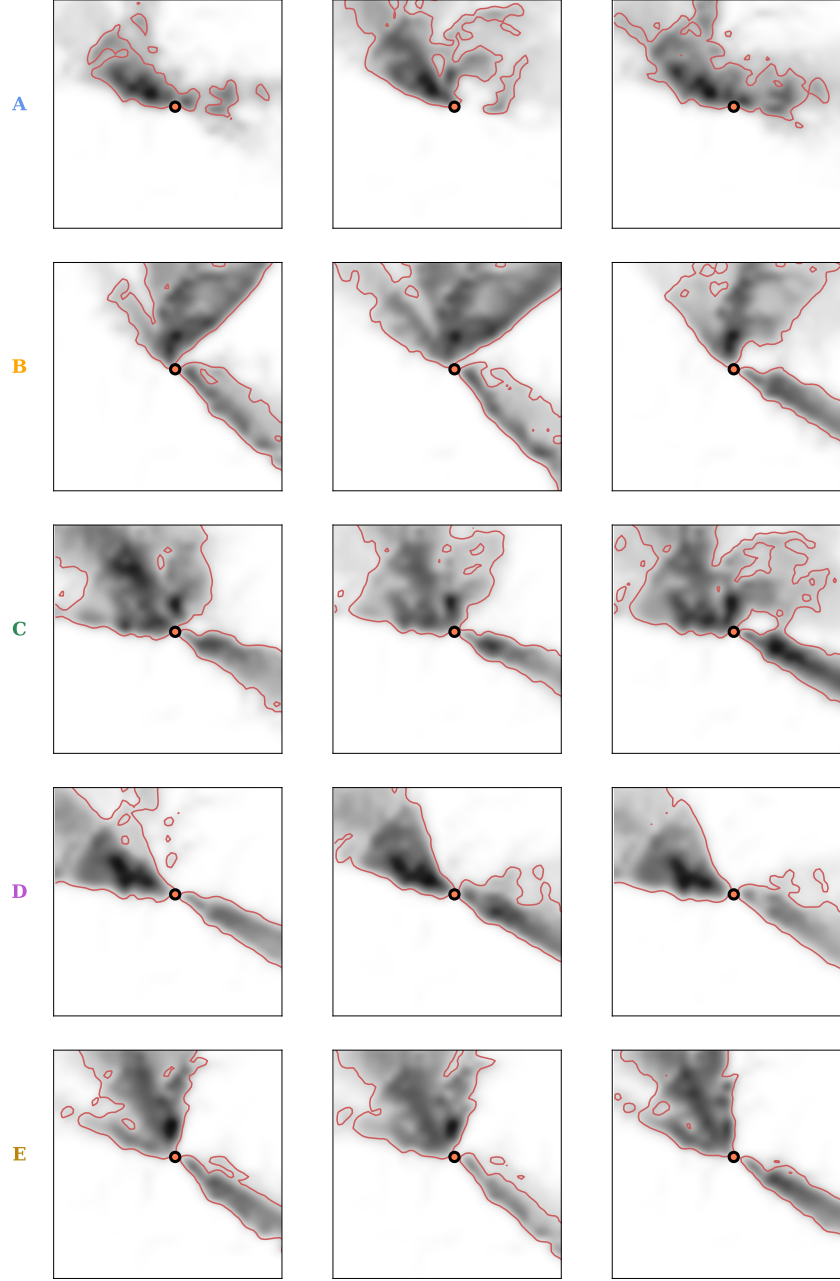


Figure 6: Five groups (in rows) of three (in column) simulated maps of maximum concentration. Grey shades denote maximum concentration (log transformed for legibility). The areas of threshold exceedance are contoured in pale red. Groups locations in the 2D AAM coordinate system are marked by large coloured circles in figure 7.

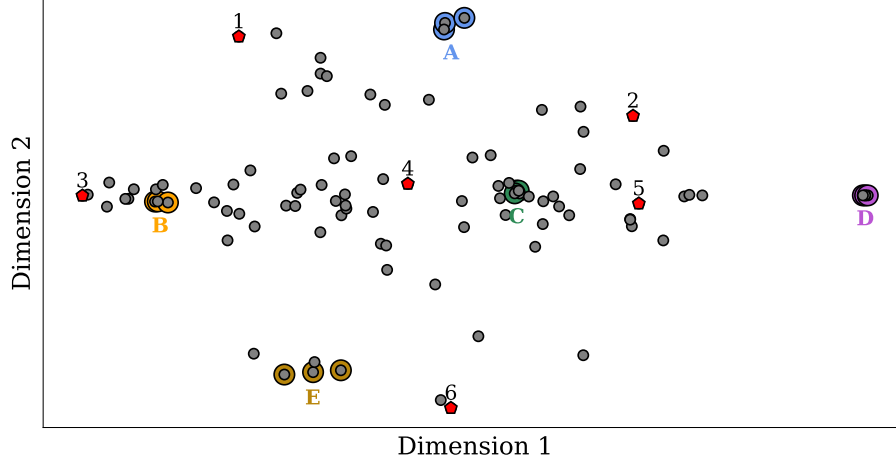


Figure 7: Scatter plot of the 100 simulations projected in the 2D AAM coordinate system. Each grey dot or red pentagon represents one simulation. Coloured circle and red pentagons highlight sets of simulations referenced in the text.

392 2D AAM coordinate system is able to capture the overall structure of the
 393 data set.

394 The left and right plots of figure 8 compare the errors induced by 2D pro-
 395 jection with PCA and AAM respectively. Each row corresponds to one sim-
 396 ulation whose location in the 2D AAM coordinate system is marked in figure
 397 7 by the red pentagon with corresponding number. Green tint indicates area
 398 where exceedance is correctly predicted after projection. Orange (respec-
 399 tively purple) tint indicates area of false positives (respectively negatives). A
 400 false positive (respectively negative) understands here as exceedance (respec-
 401 tively non exceedance) at a given location in the projected map when the
 402 threshold is not exceeded (respectively exceeded) in the original map. These
 403 examples show that AAM almost always outperform PCA: the orange and
 404 purple areas in the right column are smaller than in the left. AAM even
 405 achieve perfect reconstructions in some regions of the 2D coordinate system,
 406 for instance the row 1, 3 and 6 of figure 8. These observations also apply to
 407 the other areas of the 2D coordinate system not shown here.

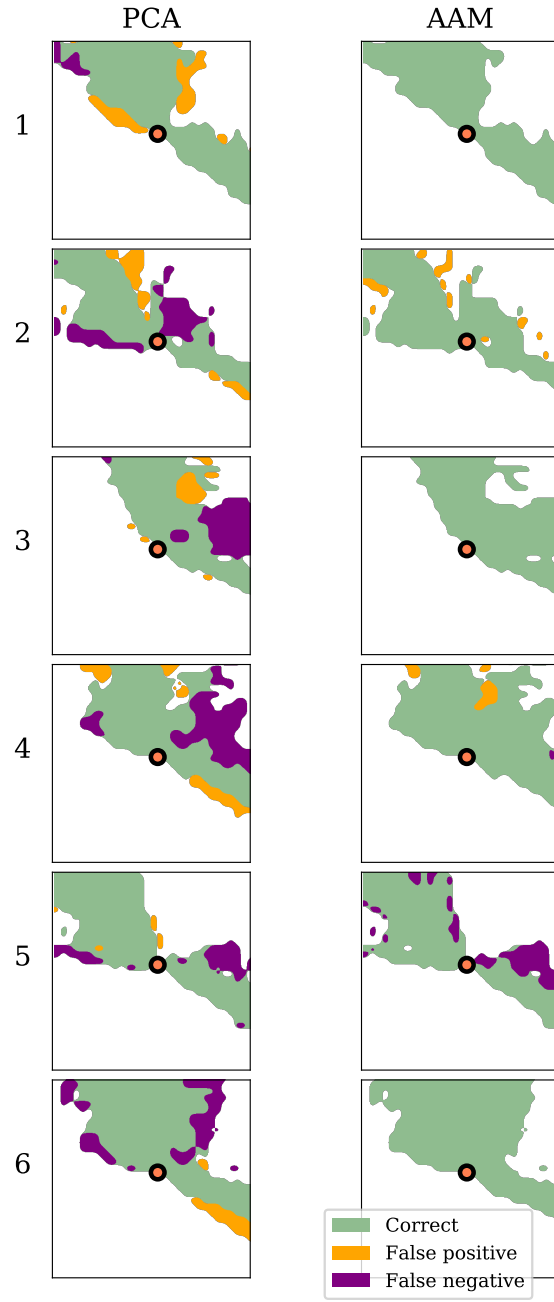


Figure 8: Comparison of errors induced by projection with PCA (left column) and AAM (right column). Each row correspond to a given original map, whose locations in the 2D AAM coordinate system are marked by numbered red pentagon in figure 7.

408 5. Conclusion

409 We proposed a generic mathematical framework aimed at modelling un-
410 certainties in 3D atmospheric dispersion simulations. It relies on stochastic
411 perturbations of both the amplitude and dynamics of the physical model in-
412 puts time series. These perturbations have a tuneable temporal structure,
413 allowing for more refined uncertainty modelling than the constant perturba-
414 tion of amplitude only that prevails in the literature.

415 To exemplify the practical use of the method, we considered a complex ac-
416 cidental situation on a rough and built-up terrain characterized by uncertain
417 release and meteorological conditions (wind direction, wind speed and pre-
418 cipitations). This realistic case study showed that the probabilistic decision
419 map obtained by uncertainty propagation can significantly impact decision.

420 The main open issue with uncertainty propagation is the limited sample
421 size than can be achieved in the short time spans characteristic of crisis con-
422 texts. In such situations, reliable exceedance probability estimates require
423 more advanced methods, such as model emulation. This raises another diffi-
424 culty, namely that those methods apply to scalar output models, not model
425 whose output is a spatial map. We argued that PCA, despite its ubiquitous
426 usage, is ill fitted for dimension reduction of such data sets, but showed that
427 AAM can overcome the shortcomings of PCA. The next step will be to lever-
428 age AAM dimension reduction to build an emulator. A possible approach is
429 to build one emulator for each AAM coordinate, for instance using Gaussian
430 process regression.

431 Contrary to situations where data is available, the structure of uncer-
432 tainty models used in crisis contexts are mostly postulated. The motivation
433 for adding a temporal structure to perturbations is to explore more exhaus-
434 tively the possible outcomes of an accident. A topic for future experimen-
435 tation would be to assess the relative influence of the parameters control-
436 ling that structure, and compare probabilistic decisions obtained with un-
437 certainty models of increasing complexity. The choice of appropriate metrics
438 for comparing decision boundaries is itself an interesting matter of investiga-
439 tion. AAM could provide additional comparison criteria based on topological
440 properties of the set of output maps.

441 References

- 442 [1] Agresti, A., Coull, B. A., may 1998. Approximate is better than “exact”
443 for interval estimation of binomial proportions. The American Statisti-

- 444 cian 52 (2), 119–126.
 445 URL <https://doi.org/10.1080/00031305.1998.10480550>
- 446 [2] Aguirre Martinez, F., Caniou, Y., Duchenne, C., Armand, P., Yalamas,
 447 T., 2016. Probabilistic assessment of danger zones associated with a
 448 hypothetical accident in a major French port using a surrogate model
 449 of CFD simulations. In: 17th International Conference on Harmonisa-
 450 tion within Atmospheric Dispersion Modelling for Regulatory Purposes,
 451 HARMO17. Budapest, Hungary.
- 452 [3] Armand, P., Brocheton, F., Poulet, D., Vendel, F., Dubourg, V.,
 453 Yalamas, T., 2014. Probabilistic safety analysis for urgent situations
 454 following the accidental release of a pollutant in the atmosphere.
 455 Atmospheric Environment 96 (0), 1–10.
 456 URL <http://www.sciencedirect.com/science/article/pii/S1352231014005433>
 457
- 458 [4] Beekmann, M., 2003. Monte carlo uncertainty analysis of a regional-
 459 scale transport chemistry model constrained by measurements from the
 460 atmospheric pollution over the paris area (ESQUIF) campaign. Journal
 461 of Geophysical Research 108 (D17).
 462 URL <https://doi.org/10.1029/2003jd003391>
- 463 [5] Bowman, V. E., Woods, D. C., jan 2016. Emulation of multivariate
 464 simulators using thin-plate splines with application to atmospheric dis-
 465 persion. SIAM/ASA Journal on Uncertainty Quantification 4 (1), 1323–
 466 1344.
 467 URL <https://doi.org/10.1137/140970148>
- 468 [6] Burgin, L., Ekström, M., Dessai, S., jan 2017. Combining dispersion
 469 modelling with synoptic patterns to understand the wind-borne trans-
 470 port into the UK of the bluetongue disease vector. International Journal
 471 of Biometeorology 61 (7), 1233–1245.
 472 URL <https://doi.org/10.1007/s00484-016-1301-1>
- 473 [7] Dabberdt, W. F., Miller, E., jan 2000. Uncertainty, ensembles and air
 474 quality dispersion modeling: applications and challenges. Atmospheric
 475 Environment 34 (27), 4667–4673.
 476 URL [https://doi.org/10.1016/s1352-2310\(00\)00141-2](https://doi.org/10.1016/s1352-2310(00)00141-2)

- [8] Duchenne, C., Armand, P., Marcilhac, M., Girard, S., Yalamas, T., 2017. A new method for assessing the uncertainty associated with 3d dispersion simulations in any variable meteorological conditions. In: 18th International Conference on Harmonisation within Atmospheric Dispersion Modelling for Regulatory Purposes, HARMO18. Bologna, Italy.
- [9] Fukunaga, K., Olsen, D. R., feb 1971. An algorithm for finding intrinsic dimensionality of data. *IEEE Transactions on Computers* C-20 (2), 176–183.
URL <https://doi.org/10.1109/t-c.1971.223208>
- [10] Galmarini, S., Bianconi, R., Klug, W., Mikkelsen, T., Addis, R., Andronopoulos, S., Astrup, P., Baklanov, A., Bartniki, J., Bartzis, J., Bellasio, R., Bompay, F., Buckley, R., Bouzom, M., Champion, H., D’Amours, R., Davakis, E., Eleveld, H., Geertsema, G., Glaab, H., Kollax, M., Ilvonen, M., Manning, A., Pechinger, U., Persson, C., Polreich, E., Potemski, S., Prodanova, M., Saltbones, J., Slaper, H., Sofiev, M., Syrakov, D., Sørensen, J., der Auwera, L., Valkama, I., Zelazny, R., sep 2004. Ensemble dispersion forecasting – part i: concept, approach and indicators. *Atmospheric Environment* 38 (28), 4607–4617.
URL <https://doi.org/10.1016/j.atmosenv.2004.05.030>
- [11] Garaud, D., Mallet, V., oct 2011. Automatic calibration of an ensemble for uncertainty estimation and probabilistic forecast: Application to air quality. *Journal of Geophysical Research* 116 (D19).
URL <https://doi.org/10.1029/2011jd015780>
- [12] Garaud, D., Mallet, V., 2012. Uncertainty estimation and decomposition based on monte carlo and multimodel photochemical simulations. Tech. Rep. 7903, Inria.
URL <https://hal.inria.fr/docs/00/67/83/06/PDF/RR-7903.pdf>
- [13] Girard, S., Iovleff, S., 2008. Auto-associative models, nonlinear principal component analysis, manifolds and projection pursuit. In: *Lecture Notes in Computational Science and Enginee*. Springer Berlin Heidelberg, pp. 202–218.
URL https://doi.org/10.1007/978-3-540-73750-6_8
- [14] Girard, S., Korsakissok, I., Mallet, V., 2014. Screening sensitivity analysis of a radionuclides atmospheric dispersion model applied to the

- 511 Fukushima disaster. *Atmospheric Environment* 95 (0), 490–500.
 512 URL [http://www.sciencedirect.com/science/article/pii/](http://www.sciencedirect.com/science/article/pii/S1352231014005317)
 513 [S1352231014005317](http://www.sciencedirect.com/science/article/pii/S1352231014005317)
- 514 [15] Girard, S., Mallet, V., Korsakissok, I., Mathieu, A., 2016. Emulation and
 515 sobol’ sensitivity analysis of an atmospheric dispersion model applied
 516 to the fukushima nuclear accident. *Journal of Geophysical Research:*
 517 *Atmospheres*.
 518 URL <http://dx.doi.org/10.1002/2015JD023993>
- 519 [16] Hanna, S., Chang, J., jan 2012. Acceptance criteria for urban dispersion
 520 model evaluation. *Meteorology and Atmospheric Physics* 116 (3-4), 133–
 521 146.
 522 URL <https://doi.org/10.1007/s00703-011-0177-1>
- 523 [17] Hanna, S. R., Paine, R., Heinold, D., Kintigh, E., Baker, D., 2007. Un-
 524 certainties in air toxics calculated by the dispersion models aermom and
 525 iscs3 in the houston ship channel area. *Journal of Applied Meteorology*
 526 *and Climatology* 46 (9), 1372–1382.
- 527 [18] Houle, M. E., 2015. Inlieriness, outlieriness, hubness and discriminabil-
 528 ity: an extreme-value-theoretic foundation. Technical Report NII-2015-
 529 002E.
- 530 [19] IMPEL, 2013. Extended release of mercaptan at a chemical plant. Tech.
 531 Rep. 43616, French Ministry for Sustainable Development.
 532 URL [https://www.aria.developpement-durable.gouv.fr/](https://www.aria.developpement-durable.gouv.fr/wp-content/files_mf/FD_43616_rouen_GB.pdf)
 533 [wp-content/files_mf/FD_43616_rouen_GB.pdf](https://www.aria.developpement-durable.gouv.fr/wp-content/files_mf/FD_43616_rouen_GB.pdf)
- 534 [20] Ismert, M., Durif, M., 2014. Accident de Lubrizol du 21 janvier 2013
 535 – couplage dispersion du nuage odorant/plaintes et appreciation des
 536 risques sanitaires associes. Tech. Rep. INERIS-DRC-13-137709-03375B,
 537 Ineris.
- 538 [21] Jolliffe, I. T., Cadima, J., mar 2016. Principal component analysis: a
 539 review and recent developments. *Philosophical Transactions of the Royal*
 540 *Society A: Mathematical, Physical and Engineering Sciences* 374 (2065),
 541 20150202.
 542 URL <https://doi.org/10.1098/rsta.2015.0202>

- 543 [22] Klampanos, I. A., Davvetas, A., Andronopoulos, S., Pappas, C.,
544 Ikonomopoulos, A., Karkaletsis, V., apr 2018. Autoencoder-driven
545 weather clustering for source estimation during nuclear events. *Envi-*
546 *ronmental Modelling & Software* 102, 84–93.
547 URL <https://doi.org/10.1016/j.envsoft.2018.01.014>
- 548 [23] Korsakissok, I., Mathieu, A., et al., 2018. Guidelines for ranking un-
549 certainties in atmospheric dispersion. Tech. Rep. Ares(2018)1172146,
550 European joint programme for the integration of radiation protection
551 research.
- 552 [24] Le, N. B. T., Mallet, V., Korsakissok, I., Mathieu, A., Perillat, R.,
553 Didier, D., Apr. 2018. Metamodeling and optimization of probabilistic
554 scores for long-range atmospheric dispersion applied to the fukushima
555 nuclear disaster. In: EGU General Assembly Conference Abstracts.
556 Vol. 20. p. 17209.
- 557 [25] Le, N. B. T., Mallet, V., Korsakissok, I., Mathieu, A., Périllat, R.,
558 01 2019. Calibration of a surrogate dispersion model applied to the
559 fukushima nuclear disaster. In: 3rd International Conference on Uncer-
560 tainty Quantification in Computational Sciences and Engineering (UN-
561 CECOMP 2019). Greece, pp. 215–228.
- 562 [26] Mallet, V., Tilloy, A., Poulet, D., Girard, S., Brocheton, F., jul 2018.
563 Meta-modeling of ADMS-urban by dimension reduction and emulation.
564 *Atmospheric Environment* 184, 37–46.
565 URL <https://doi.org/10.1016/j.atmosenv.2018.04.009>
- 566 [27] Oldrini, O., Armand, P., jan 2019. Validation and sensitivity study of
567 the PMSS modelling system for puff releases in the joint urban 2003
568 field experiment. *Boundary-Layer Meteorology* 171 (3), 513–535.
569 URL <https://doi.org/10.1007/s10546-018-00424-1>
- 570 [28] Oldrini, O., Armand, P., Duchenne, C., Olry, C., Moussafir, J., Tinarelli,
571 G., Oct 2017. Description and preliminary validation of the pmss fast re-
572 sponse parallel atmospheric flow and dispersion solver in complex built-
573 up areas. *Environmental Fluid Mechanics* 17 (5), 997–1014.
574 URL <https://doi.org/10.1007/s10652-017-9532-1>

- [29] Pearson, K., nov 1901. LIII. on lines and planes of closest fit to systems of points in space. The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science 2 (11), 559–572.
URL <https://doi.org/10.1080/14786440109462720>
- [30] Roustant, O., Ginsbourger, D., Deville, Y., 2012. DiceKriging, DiceOptim: two R packages for the analysis of computer experiments by kriging-based metamodeling and optimization. Journal of Statistical Software 51 (1), 1–55.
- [31] Skamarock, W. C., Klemp, J. B., Dudhia, J., Gill, D. O., Barker, D. M., Wang, W., Powers, J. G., 2005. A description of the advanced research wrf version 2. Tech. rep., National Center For Atmospheric Research Boulder Co Mesoscale and Microscale Meteorology Div.
- [32] Swallow, B., Rigby, M., Rougier, J., Manning, A., Lunt, M., O’Doherty, S., 2017. Parametric uncertainty in complex environmental models: a cheap emulation approach for models with high-dimensional output. arXiv preprint arXiv:1702.03696.
- [33] Tinarelli, G., Mortarini, L., Castelli, S. T., Carlino, G., Moussafir, J., Olry, C., Armand, P., Anfossi, D., mar 2013. Review and validation of MicroSpray, a lagrangian particle model of turbulent dispersion. In: Lagrangian Modeling of the Atmosphere. American Geophysical Union, pp. 311–328.
URL <https://doi.org/10.1029/2012gm001242>
- [34] Trini Castelli, S., Armand, P., Tinarelli, G., Duchenne, C., Nibart, M., nov 2018. Validation of a lagrangian particle dispersion model with wind tunnel and field experiments in urban environment. Atmospheric Environment 193, 273–289.
URL <https://doi.org/10.1016/j.atmosenv.2018.08.045>
- [35] Warner, T. T., Sheu, R.-S., Bowers, J. F., Sykes, R. I., Dodd, G. C., Henn, D. S., 2002. Ensemble simulations with coupled atmospheric dynamic and dispersion models: Illustrating uncertainties in dosage simulations. Journal of Applied Meteorology (1988-2005) 41 (5), 488–504.
URL <http://www.jstor.org/stable/26184992>

607 [36] Zhang, Y., Bocquet, M., Mallet, V., Seigneur, C., Baklanov, A., dec
608 2012. Real-time air quality forecasting, part II: State of the science,
609 current research needs, and future prospects. *Atmospheric Environment*
610 60, 656–676.
611 URL <https://doi.org/10.1016/j.atmosenv.2012.02.041>