

Accounting for representativeness errors in the inversion of atmospheric constituent emissions: application to the retrieval of regional carbon monoxide fluxes

By MOHAMMAD REZA KOOHKAN^{1,2} and MARC BOCQUET^{1,2*}, ¹*Université Paris-Est, CEREa, joint laboratory École des Ponts ParisTech and EDF R&D, Champs sur Marne, France;* ²*INRIA, Paris-Rocquencourt research center, France*

(Manuscript received 16 December 2011; in final form 4 June 2012)

ABSTRACT

A four-dimensional variational data assimilation system (4D-Var) is developed to retrieve carbon monoxide (CO) fluxes at regional scale, using an air quality network. The air quality stations that monitor CO are proximity stations located close to industrial, urban or traffic sources. The mismatch between the coarsely discretised Eulerian transport model and the observations, inferred to be mainly due to representativeness errors in this context, lead to a bias (average simulated concentrations minus observed concentrations) of the same order of magnitude as the concentrations. 4D-Var leads to a mild improvement in the bias because it does not adequately handle the representativeness issue. For this reason, a simple statistical subgrid model is introduced and is coupled to 4D-Var. In addition to CO fluxes, the optimisation seeks to jointly retrieve *influence coefficients*, which quantify each station's representativeness. The method leads to a much better representation of the CO concentration variability, with a significant improvement of statistical indicators. The resulting increase in the total inventory estimate is close to the one obtained from remote sensing data assimilation. This methodology and experiments suggest that information useful at coarse scales can be better extracted from atmospheric constituent observations strongly impacted by representativeness errors.

Keywords: inverse modelling, representativeness errors, carbon monoxide, 4D-Var

1. Introduction

In tracer transport studies, observations are infrequent in time and, for ground-measurements, sparse in space. Furthermore, they do not intrinsically carry any information about the future. That is why, complementarily, numerical models are used to assess the meteorological and chemical state of the atmosphere. In air quality modelling, input data, such as initial and boundary conditions, emission fluxes and vertical diffusion coefficients, are necessary to run proper simulations. The uncertainties of these input data and perhaps the lack of understanding of the underlying physical processes induce model errors in the simulations. To minimise them, data assimilation (DA) methods can be used. They combine observational data and information coming from chemistry and transport models and their related error statistics in

order to find the optimal values of the parameters that minimise the errors.

Four-dimensional variational DA (4D-Var) is a powerful method when it comes to constraining dynamical systems by numerous observations. In 4D-Var, all types of information mentioned above are accounted for in a two-term cost function $\mathcal{J} = \mathcal{J}_o + \mathcal{J}_b$. The first term \mathcal{J}_o is a measure of the discrepancy between the observed and simulated concentrations. The second term \mathcal{J}_b evaluates the departure of the control parameters from the first guess (background) of these parameters. By minimising the sum of these two terms, 4D-Var makes an optimal compromise while enforcing the fact that the simulated concentrations are obtained from a given numerical transport model. Iterative descent algorithms, such as conjugate gradient or quasi-Newton methods, are often used to minimise the cost function and to provide the optimal control parameters. The adjoint model is used in 4D-Var to find the gradient of the cost function with respect to these control parameters. Introducing optimal control theory ideas in

*Corresponding author.
email: bocquet@cerea.enpc.fr

geophysics, Le Dimet and Talagrand (1986) used 4D-Var to assimilate meteorological observations. Fisher and Leny (1995) used 4D-Var for the analysis of some chemically active tracer species. Lately, variational DA studies have focussed on the inverse modelling of pollutant emission fields [e.g. Elbern et al. (2007) and other references within Zhang et al. (in press)].

Focussing on carbon monoxide (CO), several modelling studies pointed out to the discrepancy between the observations and the simulated concentrations. Using the Emission Database for Global Atmospheric Research 3 (EDGAR3) inventory, before any correction, the model global run of Fortems-Cheiney et al. (2011) underestimates the CO concentrations of about 5–10% with respect to the satellite observations for January, February and March 2005. Emmons et al. (2010) compared the satellite observations to simulations of the Model for OZone And Related chemical Tracers, version 4 (MOZART-4), using the EDGAR3 inventory. Displaying a similar trend, their results exhibit an underestimation of the CO concentrations over Europe of about 10–20% for the same period.

That is why inverse modelling experiments have been carried out to update the CO flux inventories. For instance, Mulholland and Seinfeld (1995) and Saide et al. (2011) have focussed on urban scale. Yumimotoa and Uno (2006) and Kopacz et al. (2009) used 4D-Var or analytical methods to invert the emissions at regional scale. Other studies have also been performed on global scale (e.g. Pétron et al., 2002; Arellano and Hess, 2006; Stavrakou and Müller, 2006; Fortems-Cheiney et al., 2009; Kopacz et al., 2010). These studies make use of ground-based instruments that measure concentrations or they make use of satellite instruments to infer satellite-derived retrieval of CO. The former instruments are mostly used in conjunction with regional scale models whereas the latter instruments are mostly used with global scale models.

In the case of an assimilation of observations over a short period (i.e. a few hours to a few days), the parameters to be optimised are usually the initial conditions. With larger DA windows (i.e. a few days to a few months), the model is more sensitive to other parameters, such as the emissions inventory, the meteorological fields and the boundary conditions.

In most top-down (i.e. inverse modelling) studies related to the global scale, the CO emissions fluxes were found to be underestimated in the Northern Hemisphere whereas they are quite consistent with the measurements in the Southern Hemisphere (e.g. Müller and Stavrakou, 2005) or slightly overestimated (e.g. Arellano and Hess, 2006). This underestimation in the Northern Hemisphere is also found in the modelling studies (e.g. Emmons et al., 2010).

Satellite and in situ measurements require specific care when compared to transport models. The discrepancy

between the observations and the model forecast of these observations are known to be due to instrumental errors, deficiencies of the model and of the forcing fields (model error) and the *representativeness error*. The assessment of this representativeness error becomes a key issue when assimilating in situ observations, which are the focus of this paper. Indeed, the model is operative at coarser scale and, by construction, cannot simulate subgrid events. The in situ observations do capture not only the coarser scale pollutant plumes but also subgrid plumes that are not accounted for by the model. Therefore, there is a residual mismatch due to unresolved scales known as the representativeness error. In DA, it is often considered part of model error but formally ascribed to the observation error.

Due to the complexity of its estimation, an experience-based value is usually assumed for that error. This value is often chosen to be the same for all measurements. Yet that is certainly not true, because the nature of the measurements can be different (urban, rural, etc.). The maximum possible representativeness error is often chosen for all observations. Alternatively, a χ^2 criterion [used by Ménard et al. (2000) in tracer studies] can be implemented to estimate the proper magnitude of the observational errors.

In this paper, our goal is to estimate carbon monoxide surface emissions with inverse modelling, using in situ measurements from an air quality network. This network operates in France, and we wish to retrieve the emissions over France. Hence, as opposed to most of the studies mentioned earlier, the focus is on mesoscale and lower troposphere modelling. These measurements are abundant but strongly impacted by representativeness errors since many of them are influenced by nearby industrial, traffic or urban sources. Most of them aim at measuring (some of) those influential sources. To perform emission inverse modelling in this context, this lack of representativeness must be accounted for. One needs to demonstrate that observations obtained at fine scale, and strongly impacted by representativeness errors, can be assimilated with the aim of correcting a pollutant inventory defined at larger scale.

In Section 2, the atmospheric transport model (ATM) is introduced, as well as, a detailed description of the observational data. The specifications of the control space are presented. An investigation of the modelling of errors and of the uncertainties of the control parameters is also reported. In Section 3, 4D-Var is used to optimise the spatiotemporal parameters of the inventories with unsatisfactory results. Since there is a dramatic lack of representativeness of the measurements, a simple subgrid statistical model is built in order to improve the 4D-Var numerical results. The statistical model aims at taking into account the impact of close-by sources on monitoring stations. Section 4 introduces and justifies this statistical model and

its tight coupling to 4D-Var. In Section 5, the inverse modelling experiment is performed with the combination of 4D-Var and the subgrid statistical model, which will be called 4D-Var- ξ . The analysis produced by the retrieval is studied. Validations with independent observations are performed, notably using cross-validation and a long-term forecast of the CO concentrations. In Section 6, the findings of this paper are summarised. The potential and limitation of the approach are discussed.

2. Inverse modelling setup

In this section, details are given about the ingredients of the inverse modelling study: the transport model, the observations, the control variables (which are the emission parameters) and the first guess provided by the initial inventory. How to incorporate them in a 4D-Var system is described below, as well as the statistical assumptions on the errors present in the system.

2.1. Atmospheric transport model

The Eulerian chemistry and transport model Polair3D of the Polyphemus platform (Boutahar et al., 2004) is used to assess the carbon monoxide concentrations. It integrates the following transport equation:

$$\frac{\partial c}{\partial t} + \text{div}(\mathbf{u}c) = \text{div}\left(\rho \mathbf{K} \nabla \frac{c}{\rho}\right) - \Lambda c + \mathcal{R}(c) + \sigma. \quad (1)$$

Field c represents the concentration of the species, ρ the air density, \mathbf{u} the wind velocity, \mathbf{K} the turbulent diffusion tensor and σ is the volume emission term; $\text{div}(\mathbf{u}c)$, $\nabla(\rho \mathbf{K} \nabla \frac{c}{\rho})$ and Λc are the advection, diffusion and wet scavenging terms, respectively, and \mathcal{R} represents the chemical reaction term. The chemistry transport equation is completed by the initial CO concentration field c_0 at $t=0$, and the boundary condition fields $c_{\partial\Omega}$ at the boundaries $\partial\Omega$ of the domain Ω . The following condition should also be satisfied at the ground:

$$\mathbf{K} \nabla c \cdot \mathbf{n} = E - v_d c. \quad (2)$$

\mathbf{n} is the unit vector normal to the ground surface and directed upwards, v_d is the dry deposition velocity and E is the surface emission function.

All runs of the model will be performed over France. The domain extends between [41.75N, 5.25W] (the left bottom corner) and [52.75N, 12.25E] (the right top corner). The grid has the resolution of $0.25^\circ \times 0.25^\circ$. Nine vertical levels are considered from the surface up to an altitude of 2780 m. The intermediary levels are 30, 150, 350, 630, 975, 1360, 1800 and 2270 m. The meteorological fields are provided by the European Centre for Medium Range

Weather Forecasts (ECMWF). These fields have a resolution of $0.36^\circ \times 0.36^\circ$ and 60 vertical levels. The time step is 3 h. Concentrations from the global chemistry-transport model MOZART, version 2 (Horowitz et al., 2003), are used to provide boundary conditions and the initial condition. A calibration factor of 1.2 is used to correct a global underestimation of incoming carbon monoxide, following the global estimations of Emmons et al. (2010).

It has initially been examined that, within our regional, lower troposphere setup and for our timescale, carbon monoxide is barely reactive. To do so, we have compared the photochemical version of Polair3D to the tracer version (validated in Quélo et al., 2007). A small bias of $5.8 \mu\text{g m}^{-3}$ is observed between the CO concentrations with or without reactions, i.e. about 2% of the average measurements. As a consequence, neglecting the reactions, we chose to use the faster tracer version of the model.

2.2. Observations

The BDQA (Base de Données de la Qualité de l’Air, details available at <http://www.atmonet.org>) is a database listing the concentrations of several air quality pollutants over France. The (mostly hourly) collected observations are provided by 600 monitoring stations distributed all over France. For carbon monoxide, 89 stations provide hourly measurements at ground level (with an average of 75 observations per hour for the year 2005). These stations belong to one of the four different categories: industrial, traffic, urban and suburban. This gives an indication of their environment but not necessarily of their representativeness in an ATM. Larssen et al. (1999) define an area of representativeness for a station as being an area in which the concentrations do not differ from the ones measured at the station by more than a specified amount. This amount can be set to the total uncertainty of the measurement or to a value not to be exceeded in order to fulfil data quality objectives. Nappo et al. (1982) further precise that more than 90% of the concentrations measured in that area should satisfy that definition. When these conditions cannot be satisfied for a station, the latter is not deemed representative of its area.

In the case of carbon monoxide, the stations belonging to the BDQA network are far from representative as it is very difficult to determine an area of representativeness for most of them. These receptors are likely to be influenced by nearby surface fluxes (Henne et al., 2010). Background stations, far from pollution sources, are missing.

For the experiments performed in this study, 8 weeks of BDQA observations will be assimilated from 1 January 2005 to 26 February 2005, for a total of 107 914 observations, while up to more than 10 months of observations

(548 964), corresponding to the rest of the year, will be used for validation. In another experiment, about 55% of the 107 914 observations will be assimilated and the rest of the 107 914 observations will be used for validation.

The locations of the BDQA network CO monitoring stations are shown in Fig. 1.

2.3. Inventory and control variables

The first guess (background information) on the fluxes needed to perform the model runs and the inversions is provided by the anthropogenic emission from the European Monitoring and Evaluation Programme (EMEP, details can be found at <http://www.ceip.at>) inventory and the biogenic emissions of the Model of Emissions of Gases and Aerosols from Nature (MEGAN) model (Guenther et al., 2006). The EMEP inventory is modulated using hourly, weekly and monthly distribution coefficients. These coefficients are provided by the GENEMIS project (GENEMIS, 1994). The EMEP inventory has a resolution of 0.50° and the MEGAN inventory has a resolution of 0.04° . We have checked that the vegetation fire emissions over the domain defined earlier and time window of this study can be neglected.

The aim of the present study is to determine the hourly grid-size optimal sources of carbon monoxide, for both the volume source σ in eq. (1), and the emission fluxes E of eq. (2). An estimation of the number of independent control variables over a DA window of 8 weeks, a domain of 58×43 grid-cells ($0.25^\circ \times 0.25^\circ$ resolution) and six levels for the volume source, yield about 2×10^7 independent variables to retrieve. That is why we have chosen to constrain the number of degrees of freedom of control space in the following way.

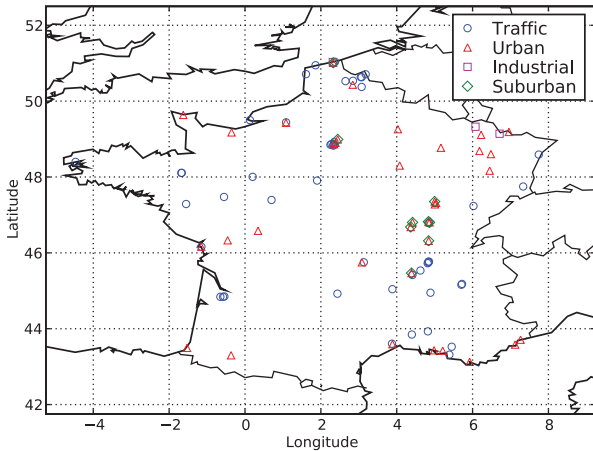


Fig. 1. The carbon monoxide monitoring stations of the BDQA network, sorted out by their official type.

The year is divided into weeks, indexed by $w = 0, \dots, N_w - 1$ where $N_w = 52$. Each week is divided into $N_h = 56$ 3-h periods, indexed by $h = 0, \dots, N_h - 1$. Each 3-h period is divided into $N_s = 3$ h, indexed by $s = 0, \dots, N_s - 1$. A grid-cell has space coordinates i, j, l (indices related to longitude, latitude and altitude, respectively) and time coordinates h, w, s [or using the global time index $k = s + N_s(h + N_h w)$]. In order to reduce the number of control variables to deal with, the discrete hourly grid-size volume sources σ and emissions E are parameterised according to

$$[\sigma]_{i,j,l,h,w,s} = [\alpha]_{i,j,h} [\sigma_b]_{i,j,l,h,w,s}, \quad (3)$$

$$[E]_{i,j,h,w,s} = [\alpha]_{i,j,h} [E_b]_{i,j,h,w,s}, \quad (4)$$

where $[\alpha]_{i,j,h}$ are the non-dimensional effective control variables corresponding to the residual degrees of freedom. They represent $58 \times 43 \times 56 = 139\,664$ scalars. The first guesses σ_b and E_b are the background sources stemming from the inventory. Let us make a remark on the temporal cycles of the inventory that are, for instance, due to vehicles traffic, urban heating, industry, etc. Because the control variables $[\alpha]_{i,j,h}$ are indexed by h , the intraweek temporal cycles will be solved for in the inverse modelling experiments. However, the longer cycles will not be solved for but are determined by the built-in cycles of the inventory: $[\sigma_b]_{i,j,l,h,w,s}$ depends on the indexes w and s . For instance, seasonal cycles of urban heating are prescribed by $[\sigma_b]_{i,j,l,h,w,s}$.

The surface emission E and volume emission σ variables have a similar local signature and would have a similar impact on a distant observation site, so that they would appear as ill-determined variables in an inverse problem. That is the reason why they were parameterised in eqs. (3) and (4) in terms of the same control vector α . It is convenient to introduce a composite emission vector e , defined in the surface layer by

$$e_{l=0} = \sigma_{l=0} + \frac{E}{\Delta}, \quad (5)$$

where Δ is the height of the surface layer. Note that this equality assumes a well-mixed surface layer. In the upper layers $l \geq 1$, it is defined by

$$e_l = \sigma_l. \quad (6)$$

In the following, the first guess about e (background) will be denoted e_b . Correspondingly, one has

$$[e_b]_{i,j,l=0,h,w,s} = [\sigma_b]_{i,j,l=0,h,w,s} + \frac{[E_b]_{i,j,h,w,s}}{\Delta} \quad \text{and} \quad (7)$$

$$[e_b]_{i,j,l \neq 0,h,w,s} = [\sigma_b]_{i,j,l \neq 0,h,w,s}.$$

As a result, eqs. (3) and (4) can be synthesised into

$$[e]_{i,j,l,h,w,s} = [\alpha]_{i,j,h} [e_b]_{i,j,l,h,w,s}. \quad (8)$$

2.4. 4D variational data assimilation

In spite of the quasi-linear physics of carbon monoxide (at these space and time scales), the computation of the Jacobian matrix is difficult to afford because of the very large set of data and control variables we intend to use. 4D-Var is meant to handle such a computational problem (Chevallier et al., 2005).

At time t_k ($k = 0, \dots, N$), the observation process is modelled with

$$\mathbf{y}_k = \mathbf{H}_k \mathbf{c}_k + \mathbf{e}_k \quad (9)$$

\mathbf{H}_k is the linear observation operator that maps the concentrations from the state space to the observation space. In this equation, $\mathbf{y}_k \in \mathbb{R}^{m_k}$ is the vector of the observed concentrations (m_k observations at time t_k), \mathbf{e}_k is the vector of observation errors at time t_k , and \mathbf{c}_k is the vector of the concentrations. The discrete form of the ATM equation, eq. (1), can be written as

$$\mathbf{c}_k = \mathbf{M}_k \mathbf{c}_{k-1} + \Delta t \mathbf{e}_k, \quad (10)$$

where \mathbf{M}_k denotes the dynamical operator of the model from t_{k-1} to t_k and Δt is the model integration time step. When t_k is only an intermediate time for model integration without observation, one has $m_k = 0$. Vector \mathbf{e}_k represents both the volume sources σ_k and the fluxes \mathbf{E}_k [see eqs. (5) and (6)].

4D-Var DA is used to invert the non-dimensional control variable vector α . The cost function to be minimised over the time-window $[t_0, t_N]$ is:

$$\begin{aligned} \mathcal{J}(\alpha) = & \frac{1}{2} \sum_{h=0}^{N_h-1} (\alpha_h - \mathbf{1})^T \mathbf{B}_{\alpha_h}^{-1} (\alpha_h - \mathbf{1}) \\ & + \frac{1}{2} \sum_{k=0}^N (\mathbf{y}_k - \mathbf{H}_k \mathbf{c}_k)^T \mathbf{R}_k^{-1} (\mathbf{y}_k - \mathbf{H}_k \mathbf{c}_k), \quad (11) \\ & + \sum_{k=1}^N \phi_k^T (\mathbf{c}_k - \mathbf{M}_k \mathbf{c}_{k-1} - \Delta t \mathbf{e}_k) \end{aligned}$$

where $k = 0, \dots, N$ is the index of integration (possibly observation) times, N_h is the number of time steps used in the time discretisation of α (in the experiments ahead $N_h = 56$), ϕ_k is a vector of Lagrange multipliers that enforces the dynamical constraint and that is called the adjoint variable, $\mathbf{R}_k = \mathbf{E}[\mathbf{e}_k(\mathbf{e}_k)^T]$ is the observation error covariance matrix, $\mathbf{B}_{\alpha_h} = \mathbf{E}[\mathbf{e}_h^b(\mathbf{e}_h^b)^T]$ is the background error covariance matrix, and $\mathbf{1}$ is the vector with entries 1. The vector α_h is the set of $[\alpha]_{i,j,h}$ for $0 \leq i \leq N_x - 1$, $0 \leq j \leq N_y - 1$ and a given h , introduced in Section 2.3. In addition, $\mathbf{e}_h^b = \alpha_h^b - \mathbf{1}$ is the background error, where α_h^b is the unknown true state of scale factors at a given h . In order to minimise the cost function \mathcal{J} with respect to α ,

with an iterative gradient-based minimiser, its gradient function can be computed as follows:

$$\begin{aligned} \nabla_{\alpha} \mathcal{J} = & \frac{\partial \mathcal{J}}{\partial \alpha} + \sum_{k=0}^{N-1} \left(\frac{\partial \mathbf{e}_k}{\partial \alpha} \right) \frac{\partial \mathcal{J}}{\partial \mathbf{e}_k} \\ = & \mathbf{B}_{\alpha}^{-1} (\alpha - \mathbf{1}) - \sum_{k=0}^{N-1} \Delta t \left(\frac{\partial \mathbf{e}_k}{\partial \alpha} \right) \phi_k. \end{aligned} \quad (12)$$

$\frac{\partial \mathbf{e}_k}{\partial \alpha}$ is a matrix, which describes the dependence of the source σ and emission \mathbf{E} as a function of the control variable vector α . Its entries can be read out from eqs. (3) and (4) and depend on \mathbf{e}_k^b .

The optimisation of eq. (11) with respect to the concentration field at time t_k gives

$$\phi_k = \mathbf{M}_{k+1}^T \phi_{k+1} + \Delta_k, \quad (13)$$

where the normalised innovation Δ_k is

$$\Delta_k = \mathbf{H}_k^T \mathbf{R}_k^{-1} (\mathbf{y}_k - \mathbf{H}_k \mathbf{c}_k). \quad (14)$$

Equation (13) is the adjoint model equation. In this equation, the boundary conditions and the final conditions are set to zero. Moreover, the bottom level condition, eq. (2), is $\mathbf{K} \nabla \phi \cdot \mathbf{n} = -v_d \phi$ (written in continuous form for the sake of simplicity).

As an approximation, the adjoint model we use is the discretisation of the continuous adjoint. This allows to use the ATM model, but propagating the concentrations backwards in time, with reversed wind fields. This approximate adjoint has been validated following Bocquet (2012), using both the so-called *duality* and *gradient* tests. For the sake of conciseness, the details are not reported here. It was checked that the errors due to the adjoint approximation are significantly smaller than the main errors' magnitude in the system.

2.5. Error modelling

In this section, we describe how the background and observation errors are statistically modelled. The background errors on the independent variables α are first related to the traditional background errors on \mathbf{e} (hence σ and \mathbf{E}). While the background error variances will be chosen a priori, the observation errors will be determined through a χ^2 diagnosis.

2.5.1 Background error covariance matrix The background error covariance matrix \mathbf{B}_{α} defines the variances–covariances between the different components of the departure of the scale factors α from $\alpha_b = \mathbf{1}$. In the inventory, anthropogenic emissions significantly dominates the biogenic emissions (1.8% of the total inventory over

France). Assuming the anthropogenic sources (such as the individual industrial sources or urban heating sources) have errors that are barely spatially correlated, the error correlation between grid-cells are taken as negligible, so that the covariance terms of that matrix are set to zero. Note that other sources of anthropogenic sources, such as traffic, might have extended correlated errors. We also neglect temporal correlations, which is a weaker assumption even though the emission are mostly anthropogenic. As a consequence of our assumptions, the prior errors are essentially represented by the variances of the prior emissions (diagonal assumption for \mathbf{B}_α).

Assuming that the emission errors are not time dependent, the variance of control variable $[\alpha]_{i,j,h}$ is

$$[\mathbf{B}_\alpha]^{i,j,h} = \frac{\sum_{w=0}^{N_w-1} \sum_{s=0}^{N_s-1} \sum_{l=0}^{N_l-1} [\mathbf{B}_e]^{i,j,l,h,w,s}}{\left(\sum_{w=0}^{N_w-1} \sum_{s=0}^{N_s-1} \sum_{l=0}^{N_l-1} [\mathbf{e}^b]_{i,j,l,h,w,s} \right)^2}, \quad (15)$$

where

$$[\mathbf{B}_e]^{i,j,l,h,w,s} = \mathbb{E} \left[\left([\mathbf{e}]_{i,j,l,h,w,s} - [\mathbf{e}^b]_{i,j,l,h,w,s} \right)^2 \right] \quad (16)$$

is the background error variance of the emission fluxes in the grid-cell of coordinates i, j, l at time h, w, s . Since the DA window of the experiments ahead is 8-week long, N_w is now set to 8.

2.5.2. Observation error covariance matrix. In eq. (9), \mathbf{e}_k includes the instrumental error and representativeness error of the observations. It is assumed that they are independent from site to site and from observation time to observation time. At this stage, the variances are assumed to be the same for all observations, which is crude since the representativeness error is expected to significantly vary between stations. Accordingly, \mathbf{R}_k is modelled as a diagonal matrix:

$$\mathbf{R}_k = r^2 \mathbf{I}_{m_k}, \quad (17)$$

where \mathbf{I}_{m_k} is the identity matrix in observation space at time t_k , and

$$r^2 = \varepsilon_{\text{repr}}^2 + \varepsilon_{\text{meas}}^2. \quad (18)$$

$\varepsilon_{\text{meas}}$ is the standard deviation of instrumental error, and $\varepsilon_{\text{repr}}$ is the standard deviation of the representativeness errors, which depends on the species, the station type and the grid size (Elbern et al., 2007).

To estimate the standard deviation parameter r , we resort to a χ^2 diagnosis ([Ménard et al., 2000; Elbern et al., 2007] for instance, in the context of atmospheric chemistry). When the statistics of the errors are consistent with the innova-

tions, then, one should expect that the average value of the cost function is equal to half of the number of assimilated observations. Accordingly, r should be chosen such that:

$$\left\{ \min_{\alpha} \mathcal{J}(\alpha) \right\}(r) \simeq \frac{m}{2}, \quad (19)$$

where $m = \sum_{k=0}^{K=N} m_k$ is the number of observations. Based on this diagnosis, an iterative process can be used to estimate r . The algorithm begins by assuming an initial value, r_0 , for r . At each iteration, r_{i+1} is computed by

$$r_{i+1}^2 = \frac{d_n^i}{m - d_s^i} r_i^2, \quad (20)$$

where d_s^i and d_n^i are twice the background part \mathcal{J}_b of the cost function and twice the observation departure part \mathcal{J}_o of the cost function, respectively, at the i th step. They respectively converge to d_s , the number of degrees of freedom for the signal (hence the s), and to d_n , the number of degrees of freedom for the noise (hence the n). The value of r is thus obtained when the sequence of r_i has converged. The method needs iterating because the minimum of the cost function does not linearly depend on r .

We note that this iterative scheme is equivalent to that of Desroziers and Ivanov (2001): eq. (20) coincides with eq. (4) of Desroziers and Ivanov (2001) when the background term is fixed. Since the method of Desroziers and Ivanov (2001) converges to one maximum of a parameter likelihood, we conclude that so does our χ^2 approach.

3. Application of 4D-Var

Following these assumptions, we perform the 4D-Var inversion of the α parameters. The assimilation window of the experiment is in the winter period, from 1 January 2005 to 26 February 2005. For comparison, a free simulation is first performed using the inventories and boundary conditions described earlier. Then, the α variables of Section 2.3 are inverted using 4D-Var.

At each grid-cell, the standard deviation of the prior error in the emission is set to 50% of the prior emission. This value is consistent with Pétron et al. (2002) and Kopacz et al. (2010). In Yumimotoa and Uno (2006), Pétron et al. (2004) and Fortems-Cheiney et al. (2009), the standard deviations are set to 100% of the prior emissions in each grid-cell, but using the EDGAR3 inventory and not over the Western Europe where the inventories are more ascertained.

An iterative test (χ^2 criterion) for the same period is applied to estimate the observational error variance. We found a standard deviation of $r \simeq 652.5 \mu\text{g m}^{-3}$ for the observational error using the χ^2 method. It is very significant since it is of the order as the average observation ($662 \mu\text{g m}^{-3}$).

A comparison of the observations with the results of the model free run, as well as a comparison to the results of the DA experiment (optimisation of α) are presented in Table 1.

The scores of this DA run show that the consistency between the analysed concentrations and the observations is low, in spite of a Pearson correlation coefficient increasing from 0.16 to 0.36. Furthermore, the reduction of the bias $\bar{O} - \bar{C}$ is unsatisfyingly small.

The total emission of the background inventory between 1st January and 26th February is 1.06 Tg. From the computation of the analysed fluxes using inverse modelling, we obtain 1.44 Tg, 36% higher than the total a priori emission. However, Fortems-Cheiney et al. (2011) estimated that value to be 17% for Western Europe, during 2005, with the reference being the EDGAR3 inventory, using biomass and anthropogenic emissions, and a spatial resolution of $2.5^\circ \times 3.5^\circ$. Kopacz et al. (2010) estimated it to be between 16% and 24% from May 2004 to April 2005. This indicates a possible over-estimation of the emission by the 4D-Var analysis. In Fig. 8 are plotted 300 h of the simulation and 4D-Var runs in the DA window, for four stations. The four corresponding profiles are too smooth to represent the peaks of the observation profile. This supports our assumption on the impact of representativeness error.

The BDQA CO network is mostly composed of proximity stations, whose observations are likely to be influenced by local sources. Therefore, the lack of consistency between the model and the observations could be explained by the direct impact of nearby pollution sources on observations. The 4D-Var analysis cannot account for the local peaks of CO concentrations since it uses a model that cannot resolve those subgrid-scale processes. However, we believe that there is some useful signal to extract from these observations. To do so, one needs to account for the subgrid processes. At least two state-of-the-art options are possible. The deterministic route consists in using explicit representations of partial information that one may have about the subgrid processes, emissions, etc. These representations are incorporated into the coarser model. This is what typically does a plume-in-grid model that uses some additional information about short-range dispersion

[e.g. Karamchandani et al. (2009) for an application to CO subgrid traffic emission]. A second route is of statistical nature. The aim is to make a statistical regression between the observations and the coarse resolution model output, which results in a fitted linear correspondence between the model to the observations. In geosciences, downscaling techniques have taken this path [e.g. Guillas et al. (2008) for an application to ozone concentrations]. In this paper, we have chosen to rely on a statistical approach to represent the subgrid effects. A deterministic modelling approach of the subgrid processes would theoretically be desirable, but it requires additional subgrid information that we do not have here, and it would be computationally more expensive.

4 Coupling 4D-Var with a subgrid statistical model

4.1. A simple subgrid statistical model

Assume that s is a continuous source field: it describes the emission at any spatial scale. Recall that \mathbf{e} is the discrete coarse-grained source that we use to drive the model. Ideally, s and \mathbf{e} should be related through a restriction, coarse-graining operator Γ , which acts as a low-pass filter, filtering out the fine details of the source:

$$\mathbf{e} = \Gamma s. \quad (21)$$

Following Bocquet et al. (2011), we can consider a prolongation operator Γ^* , which refines a coarse emission field \mathbf{e} to a continuous field s^* :

$$s^* = \Gamma^* \mathbf{e}. \quad (22)$$

There is freedom in choosing Γ^* . It could be a basic subgrid spatial interpolation operator, it could rely on additional subgrid information or it could be obtained from a Bayesian inference (Bocquet et al., 2011). For the purpose of this derivation, we do not have to specify a precise form for Γ^* . However, it is reasonable to assume $\Gamma\Gamma^* = \mathbf{I}$. Besides, $\Gamma^*\Gamma$ is a projection operator, not the identity,

Table 1. Comparison of the observations and the simulated or analysed concentrations. \bar{C} is the mean concentration, \bar{O} is the mean observation and $\text{NB} = 2(\bar{C} - \bar{O})/(\bar{C} + \bar{O})$ is the normalised bias. RMSE stands for root-mean square error. R is the Pearson correlation. FA_x is the fraction of the simulated concentrations that are within a factor x of the corresponding observations. \bar{C} , \bar{O} and the RMSE are given in $\mu\text{g m}^{-3}$

	\bar{C}	\bar{O}	NB	RMSE	R	FA_2	FA_5
Simulation (1 January–26 February 2005)	303	662	−0.74	701	0.16	0.52	0.90
Optimisation of α (4D-Var)	396	662	−0.50	633	0.36	0.59	0.92
Optimisation of ξ	615	662	−0.07	503	0.57	0.73	0.96
Coupled optimisation of α , ξ (4D-Var- ξ)	671	662	0.01	418	0.73	0.79	0.97

because of some details of the real fine scale emission field are lost in the restriction process Γ .

If \mathcal{H} is the Jacobian of a continuous multiscale hypothetical carbon monoxide model that relates s to the measurements \mathbf{y} , the vector collecting all measurements, then

$$\begin{aligned} \mathbf{y} &= \mathcal{H}s + \boldsymbol{\varepsilon} \\ &= \mathcal{H}\Gamma^*\Gamma s + \mathcal{H}(\mathbf{I} - \Gamma^*\Gamma)s + \boldsymbol{\varepsilon} \\ &= (\mathcal{H}\Gamma^*)\mathbf{e} + \mathcal{H}(\mathbf{I} - \Gamma^*\Gamma)s + \boldsymbol{\varepsilon}. \end{aligned} \quad (23)$$

Assume Γ operates the coarse-graining at the finest scale accessible by the model. Therefore, $\mathcal{H}\Gamma^*$ could be identified with the Jacobian of our Eulerian ATM. Since $\mathbf{I} - \Gamma^*\Gamma$ is a high-pass projector (it retains the short-scale fluctuations of the real emission field), $\mathcal{H}(\mathbf{I} - \Gamma^*\Gamma)s$ theoretically stands for the representativeness error (Wu et al., 2011).

Unfortunately, we do not have access to s or a multiscale model \mathcal{H} , and one needs a simple subgrid scale model to approximate $\mathcal{H}(\mathbf{I} - \Gamma^*\Gamma)s$ and close the equation. We assume this representativeness error is mostly due to subgrid/nearby sources that have a strong impact on the measurements, which are not representative of the background carbon monoxide concentration level. Another possibly significant source of error is the weakness of current vertical turbulent diffusion parameterisations. Notice that part of it may be categorised as representativeness errors when, for instance, the boundary layer height varies significantly within grid-cells.

Guided by the structure of $\mathcal{H}(\mathbf{I} - \Gamma^*\Gamma)s$, we choose to model this nearby source influence by the term

$$\xi_i \Pi_{i,k} \mathbf{e} \quad (24)$$

where ξ_i is a positive scalar attached to a station indexed by i . Similarly to $\mathcal{H}(\mathbf{I} - \Gamma^*\Gamma)s$, $\xi_i \Pi_{i,k} \mathbf{e}$ has a linear explicit dependence on the emission \mathbf{e} . The *influence coefficient* ξ_i quantifies the influence of local nearby sources onto the station. It can be interpreted as the time (given in hours in the following) required to reach a CO concentration level equivalent to the subgrid part of the measurement $[\mathbf{y} - \mathbf{H}\mathbf{c}]_{i,k}$, by emitting $\Pi_{i,k} \mathbf{e}$, which is based on the coarse-grained inventory. This influence factor is assumed constant in time and it is a priori unknown. $\Pi_{i,k}$ is an operator that linearly interpolates \mathbf{e} at the station location and at time t_k . If ξ_i is vanishing, then the representativeness of the station is deemed good. Otherwise, a significant ξ_i (a few hours and beyond) indicates a possible significant impact of nearby sources. Fig. 2 illustrates this rationale.

This term is enforced in the observation model eq. (9), which becomes, at any given time:

$$\mathbf{y} = \mathbf{H}\mathbf{c} + \boldsymbol{\xi} \cdot \Pi \mathbf{e} + \hat{\boldsymbol{\varepsilon}}, \quad (25)$$

where $\boldsymbol{\xi} \cdot \Pi \mathbf{e}$ is the vector of entries $[\boldsymbol{\xi} \cdot \Pi \mathbf{e}]_{i,k} = \xi_i \Pi_{i,k} \mathbf{e}$. The residual error $\hat{\boldsymbol{\varepsilon}}$ should statistically be smaller than $\boldsymbol{\varepsilon}$ of

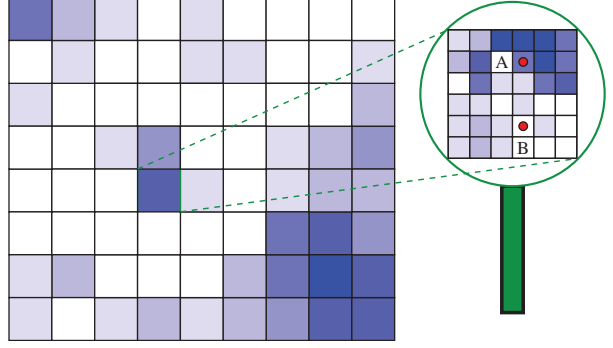


Fig. 2. Possible physical interpretation of the subgrid model. This mesh represents the CO inventory of a spatial domain. The darker the blue shade, the bigger the emission in the grid-cell. Notice the high emission zone in the south-east corner. A zoom is performed on one of the central grid-cell (see in the magnifier). Inside this grid-cell is represented a finer scale inventory inaccessible to the modeller that may represent the true multiscale inventory. Two CO monitoring stations are considered. Station A is under the direct influence of a nearby active emission zone that represents a significant contribution to the grid-cell flux. The model, operating at coarser scales, cannot scale the influence of this active zone onto station A, even though it has an estimation of its total contribution through the grid-cell total emission. Differently, station B, which is located in the same grid-cell, does not feel the active zone as much as station A. Our subgrid statistical model assumes that the influence of the active subgrid zone onto A or B has a magnitude quantified by the influence factors ξ_A and ξ_B . Obviously, in this case, one has $\xi_A \gg \xi_B$. Notice that both stations A and B are under the influence of the south-east corner of the whole domain. But this influence is meant to be represented through the Eulerian coarser ATM.

eq. (9) since part of the representativeness error should now be accounted for by the subgrid term. We denote its covariance matrix with $\hat{\mathbf{R}} = \mathbf{E}[\hat{\boldsymbol{\varepsilon}}\hat{\boldsymbol{\varepsilon}}^T]$. Under independence assumptions, the two are connected by

$$\mathbf{R} = \mathbf{E}[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T] = \boldsymbol{\xi} \cdot \Pi \mathbf{E}[\mathbf{e}\mathbf{e}^T] \Pi^T \cdot \boldsymbol{\xi}^T + \hat{\mathbf{R}}. \quad (26)$$

4.2. Coupling to the 4D-Var system

Taking into account the statistical subgrid model, the 4D-Var cost function becomes

$$\begin{aligned} \mathcal{J}(\boldsymbol{\alpha}, \boldsymbol{\xi}) &= \frac{1}{2} \sum_{h=0}^{N_h-1} (\boldsymbol{\alpha}_h - \mathbf{1})^T \mathbf{B}_{\mathbf{z}_h}^{-1} (\boldsymbol{\alpha}_h - \mathbf{1}) \\ &\quad + \frac{1}{2} \sum_{k=0}^N (\mathbf{y}_k - \mathbf{H}_k \mathbf{c}_k - \boldsymbol{\xi} \cdot \Pi \mathbf{e}_k)^T \\ &\quad \quad \times \hat{\mathbf{R}}_k^{-1} (\mathbf{y}_k - \mathbf{H}_k \mathbf{c}_k - \boldsymbol{\xi} \cdot \Pi \mathbf{e}_k) \\ &\quad + \sum_{k=1}^N \phi_k^T (\mathbf{c}_k - \mathbf{M}_k \mathbf{c}_{k-1} - \Delta t \mathbf{e}_k). \end{aligned} \quad (27)$$

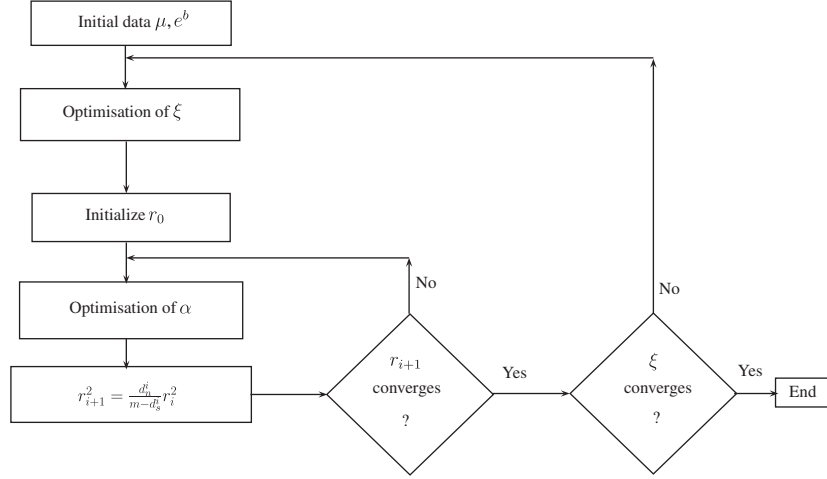


Fig. 3. Schematic of the minimisation algorithm for the 4D-Var- ξ system.

As mentioned in the previous section, if the subgrid model does account for a significant part of the representativeness error, the error covariance matrix $\hat{\mathbf{R}}_k$ should differ from \mathbf{R}_k since it accounts for the residual errors. Its magnitude will be determined by the χ^2 method.

A joint iterative optimisation of the scale factors α and the influence factor vector ξ is used to minimise the cost function. Within each iteration, ξ is obtained by a minimisation of the cost function under the constraint of positivity of the ξ_i . To perform the minimisation, one needs the gradient with respect to ξ

$$\nabla_{\xi} \mathcal{J}(\alpha, \xi) = \sum_{k=0}^N \mathbf{e}_k^T \Pi^T \hat{\mathbf{R}}_k^{-1} (\mathbf{y}_k - \mathbf{H}_k \mathbf{c}_k - \xi \cdot \Pi \mathbf{e}_k), \quad (28)$$

and the innovation vector of eq. (14) becomes

$$\Delta_k = \mathbf{H}_k^T \hat{\mathbf{R}}_k^{-1} (\mathbf{y}_k - \mathbf{H}_k \mathbf{c}_k - \xi \cdot \Pi \mathbf{e}_k). \quad (29)$$

After the ξ_i 's are optimised, the χ^2 method is used to rescale the new observational error covariance matrices $\hat{\mathbf{R}}_k = \hat{r} \mathbf{I}_{m_k}$. It is used iteratively until convergence of \hat{r} . For each cycle within this loop, the α 's are first optimised using 4D-Var for the current value of ξ and of the $\hat{\mathbf{R}}_k$. Then the $\hat{\mathbf{R}}_k$'s are updated. Fig. 3 summarises the minimisation procedure for the coupled DA system (in short 4D-Var- ξ). Note that the first step of the minimisation can begin by optimising either the influence factors ξ or the scale factor vector α . Our tests show that the final results of both minimisations are consistent. However, the former approach shows a faster convergence.

5. Application of 4D-Var- ξ

In this section, the 4D-Var- ξ system is first applied to the same setup as the 4D-Var analysis of Section 3. The resulting analysis is discussed both in terms of retrieved emission and in terms of analysed CO concentrations. Then, the system is validated with a comparison, a cross-validation and a forecast experiments.

5.1. Analysis

5.1.1. Minimisation of the cost function. Fig. 4 shows the minimisation of the cost function \mathcal{J} in the two following cases: the optimisation of the scale factor vector α (4D-Var alone) and the optimisation of α and ξ with 4D-Var- ξ .

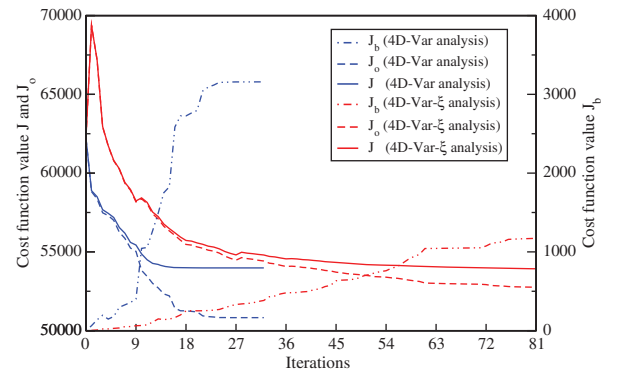


Fig. 4. Iterative decrease of the full cost function (black lines), of the background term of the cost function \mathcal{J}_b (blue lines) and of the observation departure term of the cost function \mathcal{J}_o (red lines). For the sake of clarity, the \mathcal{J}_b values are to be read on the right y-axis. Two optimisations are considered: with 4D-Var (dashed lines), and joint 4D-Var and ξ optimisation (full lines), within the assimilation window of the first 8 weeks of 2005.

In the latter case, several cycles of nine iterations each are run. In each cycle, the influence factors are first optimised and eight other iterations are used to optimise the scale factors. This cycle is repeated nine times, beyond which convergence is reached. For the first iteration of a cycle, the diagonal elements (\hat{r}) of the observational covariance matrix are diagnosed with χ^2 . This may lead to a temporary increase of the cost function value as seen in Fig. 4. In both cases, the cost function \mathcal{J} consistently converges to half of the observation numbers (that is, $m/2 = 53,957$). The values of the observation and background terms of the cost function, \mathcal{J}_o and \mathcal{J}_b respectively, have also been plotted (cf. Fig. 4).

The \mathcal{J}_o of 4D-Var- ξ converges to a higher value than the \mathcal{J}_o of 4D-Var because the coupled scheme is able to identify a higher fraction of the degrees of freedom as noise (representativeness errors). The \mathcal{J}_b of 4D-Var- ξ converges to a smaller value than the \mathcal{J}_b of 4D-Var because the coupled scheme recognises that the degrees of freedom for the signal present in the observations are significantly less important than what 4D-Var would assume. Specifically, the number of degrees of freedom for the signal is $d_s = 6316$ with 4D-Var, whereas it is $d_s = 2367$ with 4D-Var- ξ . They stand for about 2% of the information load of the in situ observations. This shows that ignoring the representativeness issue leads to a severe overestimation of the information content of the dataset. The standard deviation of the residual diagnosed observation error that was $r \simeq 652.5 \mu\text{g m}^{-3}$ without the implementation of the subgrid scheme is now $\hat{r} \simeq 422 \mu\text{g m}^{-3}$.

5.1.2. Scores. Statistical indicators are computed for the output of an 8-week experiment using the 4D-Var- ξ scheme. They are reported in Table 1 (joint optimisation of ξ and α). A significantly better agreement is obtained between the analysis and the observations. The large underestimation of the CO concentrations (see the means in Table 1) is significantly reduced: the normalised bias is as small as 1.4%. The total emission is diagnosed to be 1.16 Tg. This is an inventory increase of about 9%, which is rather consistent with studies performed over Western Europe using remote sensing. In addition to the bias reduction, it also leads to an increase of the Pearson correlation coefficient up to 0.73. The optimisation of the influence coefficients, using the a priori fluxes, leads to decrease the root mean square error (RMSE) from $701 \mu\text{g m}^{-3}$ to $503 \mu\text{g m}^{-3}$. The emission optimisation decreases this number down to $418 \mu\text{g m}^{-3}$. The impact of the subgrid model on the RMSE is consistent with the predominance of the local sources on the observations.

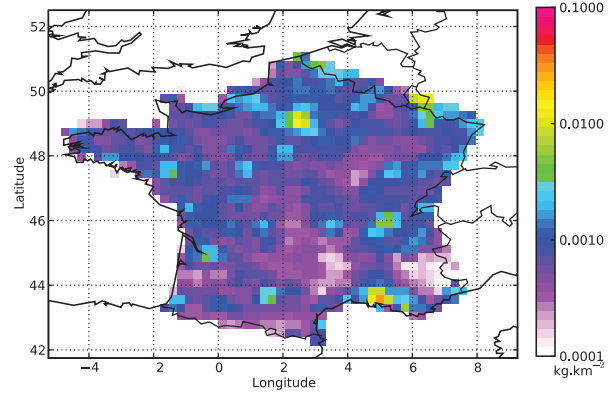


Fig. 5. Time-integrated spatial distribution of the carbon monoxide EMEP+MEGAN inventory over the first 8 weeks of 2005.

5.1.3. Spatial distribution of the retrieval. The values of the scale factors α of the 4D-Var- ξ system range between 0.01 and 19.5, with an average value of 1, showing that some important correction can be made to the inventory. Fig. 5 displays the carbon monoxide EMEP+MEGAN inventory (the first guess) integrated over the first 8 weeks of 2005, for each grid-cell. Fig. 6 displays the ratio of time-integrated retrievals to the time-integrated EMEP+MEGAN inventory, for each grid-cell. Fig. 6a displays the retrieval obtained using 4D-Var, whereas Fig. 6b displays the retrieval obtained using 4D-Var- ξ . 4D-Var- ξ shows a much less pronounced correction than the 4D-Var retrieval, which is consistent with the findings from the statistics discussed in the previous section. The joint inverse modelling retrieval suggests an increase of the emissions in the South of Paris area, Lyons, La Rochelle, Lille and in the Mediterranean coast of France, pointing to an underestimation of the inventory. It suggests a decrease of the emissions in the area of Dunkerque, Metz and North of Paris, pointing to an overestimation of the inventory.

5.1.4. Results: scatterplots. In Fig. 7a, a scatterplot compares the observations to the concentrations simulated by the model using the a priori emissions. It is clearly impacted by the representativeness errors, since the variability of the observations is much stronger than that of the simulated concentrations. In Fig. 7b, a second scatterplot compares the observations to the ATM concentrations using the a posteriori emissions from 4D-Var. Even though 4D-Var corrects the shape of the scatterplot, it is still highly impacted by representativeness errors. Fig. 7c is a scatterplot of the observations versus the concentrations diagnosed by the 4D-Var- ξ system. The representativeness errors have been significantly reduced. However, there is still a residual impact for the smallest observations.

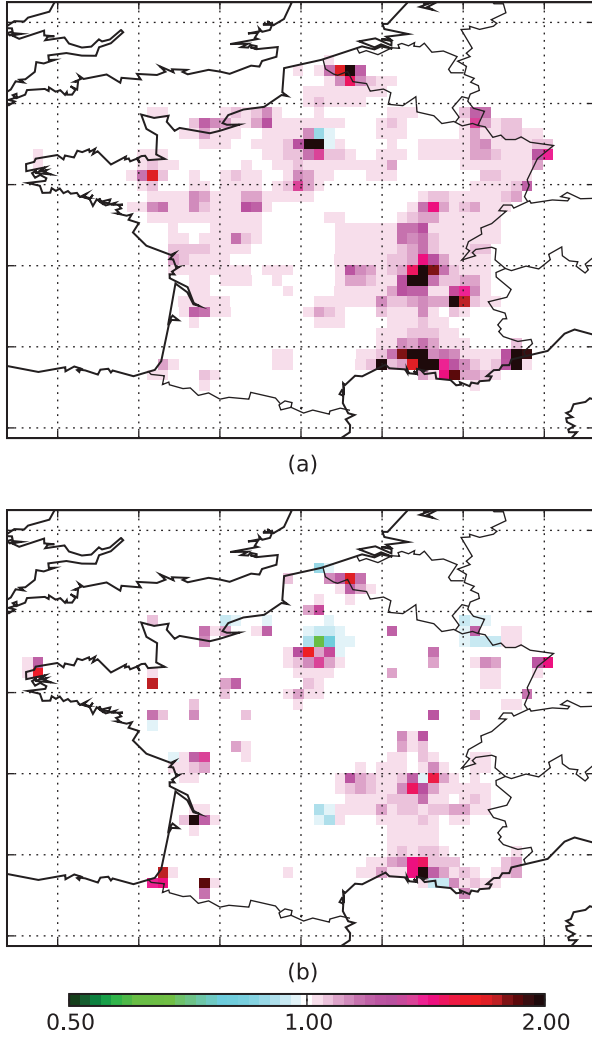


Fig. 6. Ratio of the time-integrated CO flux retrieval to the EMEP+MEGAN time-integrated CO flux for each grid-cell, in the 4D-Var case (a) and in the joint 4D-Var and subgrid model case (b).

This may be due to situations where carbon monoxide emitted locally is not advected nearby monitoring station i , whereas ξ_i may be significant because of the impact of the local source when the winds are blowing in the direction of the instrument. Indeed, our simple statistical model cannot account for the changes in the local micrometeorology, only for its indirect impact.

5.1.5. On-site profiles. Here, the focus is on the analysis at individual stations. The values of the station-dependent influence factors ξ_i range between 0 and 97.5 h, with a median value of 5.9 h and a mean value of 11.3 h.

In Fig. 8, four different time series of concentrations are displayed for four different stations: the observations, the

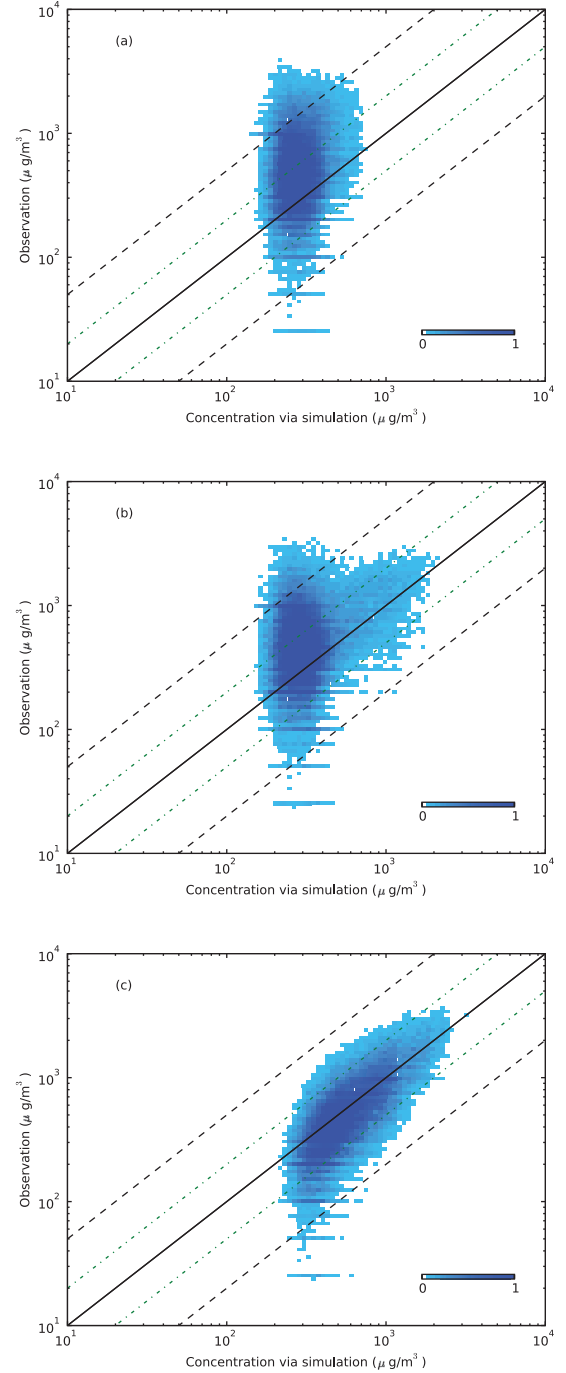


Fig. 7. Scatterplot during 8 weeks: (a) comparison between the concentrations via the model and the observations, (b) comparison between the concentrations via the model using the a posteriori emissions retrieved from 4D-Var and the observations, (c) comparison between the concentrations diagnosed by the 4D-Var- ξ system and the observations. The colour bars show the correspondence between the blue shade and the density of points of the scatterplot. This density has been normalised so that its maximum is 1. Dashed lines are the FA₅ dividing lines, and dashed-dotted lines are the FA₂ dividing lines.

concentrations simulated with the a priori emissions, the concentrations obtained from 4D-Var and 4D-Var- ξ concentrations. The traffic station of Lille Pasteur, can be cited as an example of small influence factor value with $\xi_i = 0.6$ h. In that station, the simulation concentrations are in quite good agreement with the observations. The correlation between the observations and the simulated concentrations reaches 0.49. It is 0.74 for the 4D-Var- ξ results. At the station Paris, boulevard périphérique

Auteuil (suburban), for which ξ_i is of 2.7 h, the correlation increases from 0.29 up to 0.77. Orléans Gambetta (traffic zone) station can be cited as an example with a moderate influence factor value of $\xi_i = 11.9$ h. At this station, the Pearson correlation coefficient increases from 0.11 to 0.67 when using the 4D-Var- ξ system. The dependence of the observations and the local emissions is clearly shown in Figure 8c. The model simulation gives a smooth curve, whereas the observations are highly fluctuating. The 4D-Var

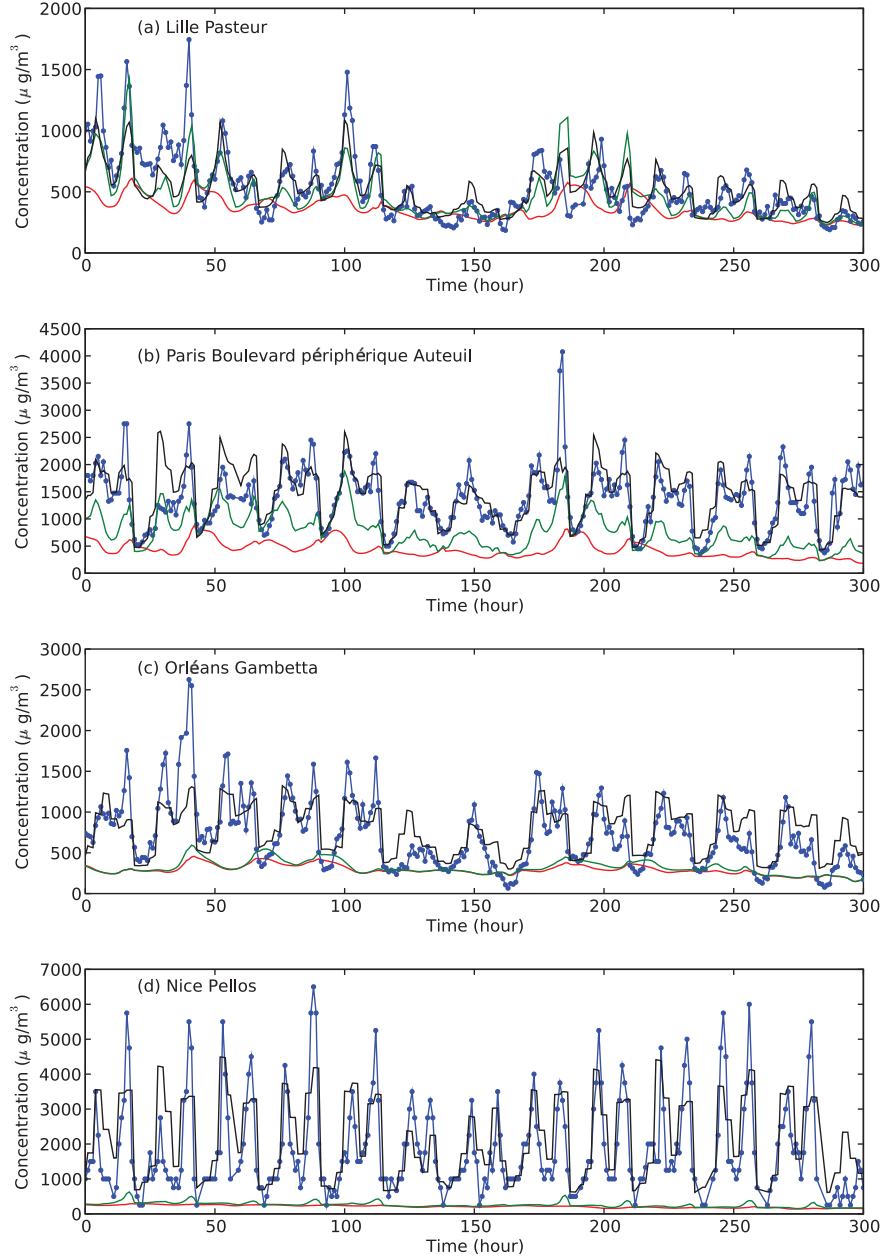


Fig. 8. Time series of CO concentrations for the first 300 h of 2005, at four stations: observations (blue), simulation using the prior emissions (red), simulation using the posterior emissions of data assimilation (green) and simulation using the posterior emissions of 4D-Var- ξ (black) with adjusted observations using the statistical subgrid model.

system is able to anticipate the trend of the concentrations but cannot predict the peaks. Furthermore, it overestimates the inventory by trying to adjust to the peaks.

Figure 8d shows the concentrations in Nice Pellos (urban station) with a high influence factor value of $\xi_i = 45.8$ h. The results of 4D-Var- ξ are in good agreement with the observations whereas neither the simulation nor 4D-Var is able to match the observations. The correlation value is significantly increased from 0.32 to 0.68. It is also clear that, although 4D-Var- ξ is able to account for a substantial part of the peaks, it underestimates their maxima and overestimates the minima, which may be due to residual representativeness error.

5.2. Validation

A direct and reliable validation of a spatial emission inventory is currently out of reach for most pollutants [see the in-depth discussion of Vestreng et al. (2007) about SO_2]. It is only possible to compare with another independent estimation (top-down or bottom-up), which, as a relative comparison approach, may not be as satisfying as a straight comparison to observations. Local flux measurements are possible (e.g. for CO_2) in some media but these are sparse and cannot fully validate a spatial inventory. Therefore, a CO emission inventory can only be indirectly validated. For instance, one can compare the CO concentrations simulated with the inventory to real measurements.

We shall first compare the total emitted carbon monoxide to an independent bottom-up inventory over France. We will then compare simulated concentrations obtained with an inventory retrieved from a training network, on a distinct validation network. Finally, after an assimilation period of 8 weeks, we shall make a 10-month CO concentration forecast. The forecasted concentrations will be compared to independent observations (that have not been assimilated).

5.2.1. Global comparison with the CITEPA inventory. The total retrieved CO emitted mass from 4D-Var- ξ is compared to the inventory of the Centre Interprofessionnel Technique d'Etudes de la Pollution Atmosphérique (CITEPA, http://www.citepa.org/emissions/nationale/Aep/aep_co.htm). According to CITEPA, the total French inventory for 2005 is 5.3 Tg. We have inferred the total emitted mass for the first 8 weeks of 2005 using the weekly and the monthly coefficients of GENEMIS for each of the 11 sectors of the SNAP nomenclature of emitting activities. The contribution of each SNAP sector to the total emission is estimated following EMEP distribution for this year. Following this rationale, the total CO emitted mass of the CITEPA inventory is found to be 1.15 Tg between

1st January and 26th February. This value is very close to 1.16 Tg obtained with 4D-Var- ξ .

5.2.2. Cross-validation experiment. Forty-nine BDQA stations have been randomly selected as a training network. Inverse modelling will be performed using the CO observations of this subnetwork for the first 8 weeks of 2005. The rest of the stations of the BDQA network forms a 40-station validation network. The observations of these stations will be compared to the simulated CO concentrations obtained using the retrieved emission field inferred from the training set. The partition between the BDQA stations is displayed in Fig. 9.

Three simulations for validation are performed: a simulation using the EMEP+MEGAN background inventory; a simulation using the emissions retrieved with 4D-Var; and a simulation using the emissions retrieved with 4D-Var- ξ . In addition to these three simulations, we shall use the influence coefficients ξ_i attached to the stations of the validation network to correct the concentrations, using the background emissions, the 4D-Var retrieved emissions and the 4D-Var- ξ retrieved emissions. Even though these 40 factors have been inferred (in the previous section) using observations of the full network, we believe they are intrinsic to the stations. Inferring them from a different (sufficiently large) observation set would yield close values. We have checked this by comparing the ξ_i of the training network obtained from a 89-station (full network) optimisation, with the ξ_i of the training network obtained from a 49-station (training network) optimisation. The results, that are reported in a scatterplot Fig. 10, confirm that the

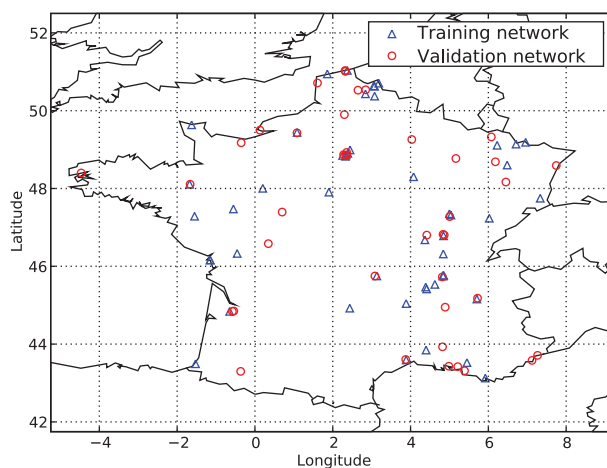


Fig. 9. The training (triangle) and validation (circle) subnetworks that partition the BDQA stations measuring carbon monoxide. This partition is randomly generated for the cross-validation experiment.

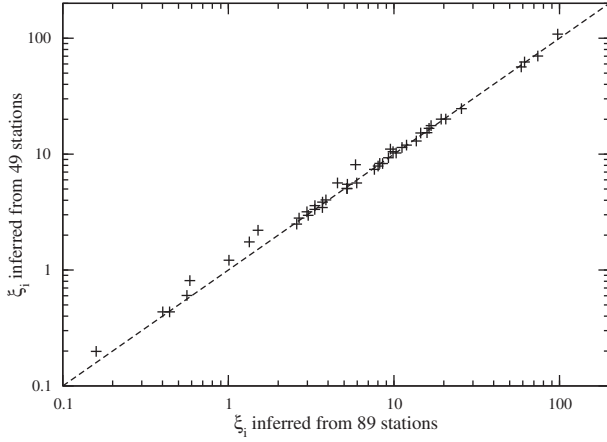


Fig. 10. Scatterplot of the 49 ξ_i of the training network inferred from either the training network or the full network (89 stations). Four $\xi_i=0$ crosses are missing. In the four cases, they were concordantly diagnosed to be 0 by the two inferences.

values are close and support that they are intrinsic to each station.

The statistical scores, as well as the total emitted mass, for these six validation experiments are reported in Table 2.

Firstly, 4D-Var- ξ without correction at the validation stations performs poorly, with scores of the same order as 4D-Var. This is to be expected since 4D-Var- ξ is meant to be used in conjunction with the ξ coefficients, which is not the case for this experiment. Secondly, 4D-Var yields sensibly better scores than 4D-Var- ξ . This is due to the excessive correction of 4D-Var that wrongly takes the CO peaks as a systematic bias. As should be, this bias correction equally applies to the validation set, leading to slightly better scores than 4D-Var- ξ but for the wrong reasons.

Applying the ξ_i coefficients of the validation stations to the concentrations obtained with the first guess emissions considerably reduces the bias and improves all the other statistical indicators as compared to the reference simulation. Applying the ξ_i coefficients of the validation stations to the concentrations obtained with the 4D-Var retrieved

emissions leads to a very large positive bias. Even though the approach is by construction inconsistent, it yields significantly better scores as compared to using the 4D-Var retrieval without corrections on the validation stations. Lastly, the ξ_i coefficients of the validation stations are used in conjunction with the 4D-Var- ξ retrieved emission field. This leads to much higher scores than the other experiments. These indicators are consistent with the scores obtained using the full network data (in Table 1).

It is remarkable that the total retrieved mass of this last experiment, 1.14 Tg, is consistent with that obtained by 4D-Var- ξ using all stations, that is, 1.16 Tg. A convincing validation of such a retrieval methodology would require such a consistency. The same is not true for 4D-Var with 1.25 Tg obtained using the training subnetwork and 1.44 Tg using the full network, pointing to the inconsistency of the method that does not properly account for the representativeness errors.

5.2.3. Forecast experiments. A validation forecast is performed over the year 2005. This second indirect validation is demanding since no new observation are assimilated over a 10-month period. That is why, in atmospheric chemistry/air quality, a forecast is often considered a more stringent validation test (Zhang et al., in press). However, our validation by a forecast has a limitation due to the statistical subgrid model. It is meant to efficiently apply to the observational network employed in the initial assimilation time-window. Notice that this limitation is inherent to any forecasting system making use of some form of statistical adaptation.

Four runs are considered. They all use the ECMWF meteorological fields and the MOZART, version 2, output for the initial and boundary conditions. The first run is a direct simulation over 2005 that is driven by the EMEP+MEGAN inventory. The second one is a direct run from 26th February to 31st December, but using the optimal α obtained from the 4D-Var analysis from 1st January to 25th February and eq. (8) to generate the inventory. The third one is a direct run from 26th February

Table 2. Comparison of the observations and the forecasted concentrations on the validation network for the first 8 weeks of 2005. The statistical indicators are described in Table 1. Additionally, the total retrieved emitted mass is given (in Tg). The corresponding value for the retrieved mass using the full network is recalled in parenthesis

Used inventory	\bar{C}	\bar{O}	NB	RMSE	R	FA ₂	FA ₅	Total mass
Background	296	697	-0.81	771	0.16	0.51	0.88	1.06 (1.06)
4D-Var	357	697	-0.65	726	0.28	0.57	0.89	1.25 (1.44)
4D-Var- ξ	310	697	-0.77	758	0.22	0.52	0.89	1.14 (1.16)
Background + climatological ξ	644	697	-0.08	538	0.60	0.73	0.96	1.06 (1.06)
4D-Var + climatological ξ	968	697	0.33	1216	0.40	0.67	0.94	1.25 (1.44)
4D-Var- ξ + climatological ξ	674	697	-0.03	514	0.64	0.75	0.96	1.14 (1.16)

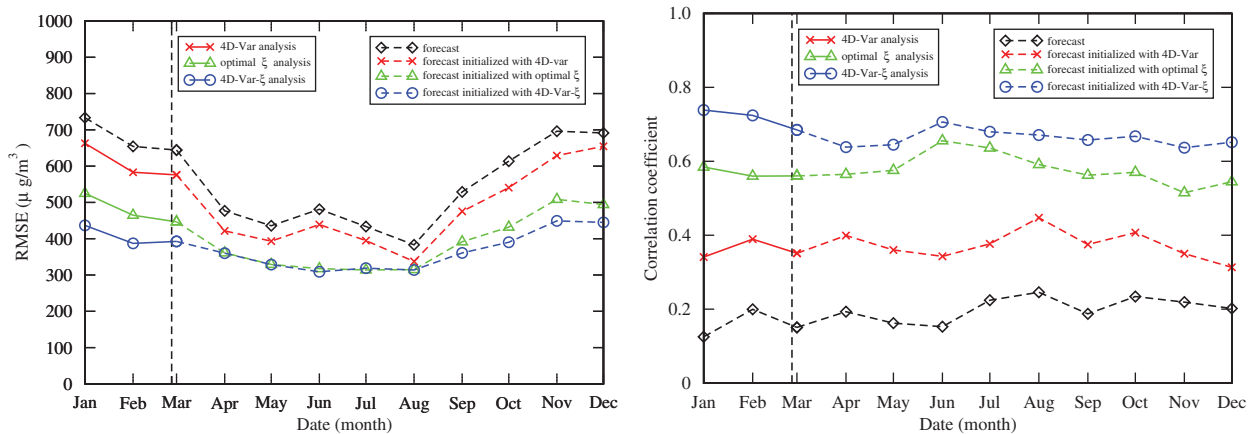


Fig. 11. Monthly RMSE (left panel) and Pearson correlation (right panel) of four runs: a pure forecast, a 10-month forecast initialised by an 8-week 4D-Var assimilation, a 10-month forecast initialised by an 8-week window where the ξ 's are optimised and a 10-month forecast initialised with an 8-week joint 4D-Var and ξ optimisation. The vertical dashed line indicates the end of the assimilation window and the start of the forecasts.

to 31st December, using the EMEP+MEGAN inventory but using the optimal ξ obtained from an optimisation over ξ of the total cost function from 1st January to 25th February. The fourth one is a direct run from 26th February to 31st December but using the optimal α and ξ parameters obtained from the 4D-Var- ξ analysis from 1st January to 25th February and eq. (8) to generate the inventory. None of the observations from 26th February to 31st December are assimilated. They are exclusively used for validation.

Such forecast requires a forecast of the emissions. The parameterisation of the emission by the α allow us to do so. In particular, some of the temporal (but not spatial) seasonal variability is implicitly accounted for thanks to the GENEMIS temporal modulation present in the first guess \mathbf{e}_b .

Firstly, we have focussed on the first month forecast, from 26th February to 26th March, where one can assume that the winter emission trend endures. The results are in very good agreement with the observations. For the forecast period, the correlation coefficient between the observations and 4D-Var- ξ increases from 0.13 to 0.68. The RMSE is improved by about 40% during the analysis period. Almost 68% of that improvement is due to the optimisation of the influence factors ξ_i .

Secondly, we have extended the forecast period, from 26th February to 31st December across seasons. The monthly results for the RMSE and the correlation coefficients, over the year 2005, are presented in Fig. 11. Using 4D-Var- ξ , the RMSE decreases by $282 \mu\text{g m}^{-3}$ within the analysis period, 1st January to 26th February (left side of the vertical dashed line). It decreases by $172 \mu\text{g m}^{-3}$ during the forecast period, from 26th February to 31st December (right side of the vertical dashed line). The improvement is

remarkably persistent during the whole 10-month forecast period. It shows that choosing α and ξ as control vectors has a good prognostic value. In spring and summer, the RMSE decreases for all four experiments. This can be due to the decrease of urban heating during that period, which is accounted for in the cycles of the inventory but which reduces a source of uncertainty. It can also be seen that the RMSE gain in the spring and summer is essentially due to the subgrid model identification, and not the emission estimation, since 4D-Var- ξ and the optimal- ξ forecast yield the same RMSE. Unsurprisingly, this means that the emission retrieval carried out over two winter months are not optimal for the spring and summer months. Another possible explanation is the emergence of new source of errors in the spring–summer time, such as the higher OH concentration that leads to a higher reactivity of CO or a stronger turbulent mixing in the boundary layer. However, this should be balanced by a persistent gain in the spring–summer period of the correlation due to the emission retrieval.

6. Conclusion

In this article, a 4D-Var DA system was developed to estimate carbon monoxide fluxes at regional scale. An approximate adjoint of the Polair3D model has been built and validated for this 4D-Var system. A study over France, at a resolution of $0.25^\circ \times 0.25^\circ$, is conducted. We used the in situ observations of the BDQA database that includes the observations from industrial, traffic, urban and suburban stations. They are strongly impacted by local sources that the stations are meant to monitor. Hence, although the number of observations is very significant, their information load is impacted by large representativeness

errors. The Pearson correlation coefficient between the simulated concentrations and the observations is computed to be 0.16. A first 4D-Var inversion of the CO fluxes leads to a mild improvement of the skill. The Pearson correlation climbs to 0.36. However looking at stations profile, it is clear that the representativeness errors are not accounted for, since the analysis from 4D-Var cannot reproduce the intense CO peaks. Besides, it leads to an artificially large increase of the retrieved emissions.

Therefore, a simple model is developed to statistically represent the subgrid effects of nearby sources. A coefficient attached to each station is used to estimate this influence. The 4D-Var system is coupled to this subgrid model and the fluxes are determined altogether with the influence coefficients. The correlation coefficient reaches 0.73, while the bias between the observations and the analysed concentrations is considerably reduced. The net increase of the CO inventory is estimated to be 9%, consistent with other top-down approaches using satellite data. Cross-validation experiments using a training subnetwork and a validation subnetwork demonstrates the consistency of the inventory estimation, whereas, in this context, the traditional 4D-Var does not deliver consistent estimations with different training subnetworks. Forecast experiments with the analysed coefficients and fluxes over 10 months, after an assimilation window of 8 weeks, show remarkably persistent scores throughout the year. This emphasises the relevance of the choice of ξ and α as joint control parameter vectors of the 4D-Var- ξ analysis.

We believe that this methodology and experiment show that, in this context, it is possible to extract relevant information from observations strongly impacted by representativeness errors. One limitation that is inherent to the statistical adaptation component of the system is that it is meant to be used on a given monitoring network. A validation forecast can safely be made to additional stations, but statistical adaptation cannot be performed to these stations, if the related influence factor ξ_i were not previously estimated.

To improve the present statistical subgrid model, which uses the influence factors to estimate the immediate impact of the emissions on the observations, a more comprehensive statistical subgrid model could be used. For instance, that model could include the effects of the wind direction, deposition parameters, etc., which are used or diagnosed in the coarse resolution model. Computationally, it would not be as cheap as the subgrid model used here.

Beyond the carbon monoxide context of this paper, we believe that the integration of the simple statistical subgrid scale into a 4D-Var can be generalised to pollutants whose observations could highly be impacted by representativeness errors.

7. Acknowledgements

This article is a contribution to the MSDAG project supported by the *Agence Nationale de la Recherche*, grant ANR-08-SYSC-014, and a contribution to the INSU/LEFE ADOMOCA-2 project. We are grateful to two anonymous reviewers for their useful comments and suggestions at the origin of several developments in the article.

References

- Arellano, A. F. and Hess, P. G. 2006. Sensitivity of top-down estimates of CO sources to GCTM transport. *Geophys. Res. Lett.* **33**, L21807.
- Bocquet, M. 2012. Parameter field estimation for atmospheric dispersion: application to the Chernobyl accident using 4D-Var. *Q. J. Roy. Meteor. Soc.* **138**, 664–681. DOI: 10.1002/qj.961.
- Bocquet, M., Wu, L. and Chevallier, F. 2011. Bayesian design of control space for optimal assimilation of observations. I: consistent multiscale formalism. *Q. J. Roy. Meteor. Soc.* **137**, 1340–1356.
- Boutahar, J., Lacour, S., Mallet, V., Musson-Genon, L., Quélo, D. and co-authors. 2004. Development and validation of a fully modular platform for the numerical modeling of air pollution: POLAIR. *Int. J. Environ. Pollut.* **22**, 17–28.
- Chevallier, F., Fisher, M., Peylin, P., Serrar, A., Bousquet, P. and co-authors. 2005. Inferring CO₂ sources and sinks from satellite observations: method and application to TOVS data. *J. Geophys. Res.* **110**, D24309.
- Desroziers, G. and Ivanov, S. 2001. Diagnosis and adaptive tuning of observation-error parameters in a variational assimilation. *Q. J. Roy. Meteor. Soc.* **127**, 1433–1452.
- Elbern, H., Strunk, A., Schmidt, H. and Talagrand, O. 2007. Emission rate and chemical state estimation by 4-dimensional variational inversion. *Atmos. Chem. Phys.* **7**, 3749–3769.
- Emmons, L. K., Walters, S., Hess, P. G., Lamarque, J.-F., Pfister, G. G. and co-authors. 2010. Description and evaluation of the model for Ozone and related chemical tracers, version 4 (MOZART-4). *Geosci. Model Dev.* **3**, 43–67.
- Fisher, M. and Leny, D. J. 1995. Lagrangian four-dimensional variational data assimilation of chemical species. *Q. J. Roy. Meteor. Soc.* **121**, 1681–1704.
- Fortems-Cheiney, A., Chevallier, F., Pison, I., Bousquet, F., Carouge, C. and co-authors. 2009. On the capability of IASI measurements to inform about CO surface emissions. *Atmos. Chem. Phys.* **9**, 8735–8743.
- Fortems-Cheiney, A., Chevallier, F., Pison, I., Bousquet, P., Szopa, S. and co-authors. 2011. Ten years of CO emissions as seen from measurements of pollution in the troposphere (MOPITT). *J. Geophys. Res.* **116**, D05304.
- GENEMIS. 1994. *Generation of European Emission Data for Episodes (GENEMIS) Project*. Technical Report, EUROTRAC Annual Report 1993, Garmisch-Partenkirchen, Germany.
- Guenther, A., Karl, T., Harley, P., Wiedinmyer, C., Palmer, P. and co-authors. 2006. Estimates of global terrestrial isoprene emissions using MEGAN (Model of Emissions of Gases and Aerosols from Nature). *Atmos. Chem. Phys.* **6**, 3181–3210.

- Guillas, S., Bao, J., Choi, Y. and Wang, Y. 2008. Statistical correction and downscaling of chemical transport model ozone forecasts over Atlanta. *Atmos. Environ.*, **42**, 1338–1348.
- Henne, S., Brunner, D., Folini, D., Solberg, S., Klausen, J. and co-authors. 2010. Assessment of parameters describing representativeness of air quality in-situ measurement sites. *Atmos. Chem. Phys.* **10**, 3561–3581.
- Horowitz, L. W., Walters, S., Mauzerall, D. L., Emmons, L. K., Rasch, P. J. and co-authors. 2003. A global simulation of tropospheric ozone and related tracers: description and evaluation of MOZART, version 2. *J. Geophys. Res.* **108**, D24.
- Karamchandani, P., Lohman, K. and Seigneur, C. 2009. Using a sub-grid scale modeling approach to simulate the transport and fate of toxic air pollutants. *Environ. Fluid. Mech.* **9**, 59–71.
- Kopacz, M., Jacob, D. J., Fisher, J. A., Logan, J. A., Zhang, L. and co-authors. 2010. Global estimates of CO sources with high resolution by adjoint inversion of multiple satellite datasets MOPITT, AIRS, SCIAMACHY, TES. *Atmos. Chem. Phys.* **10**, 855–876.
- Kopacz, M., Jacob, D. J., Henze, D. K., Heald, C. L., Streets, D. G. and co-authors. 2009. A comparison of analytical and adjoint Bayesian inversion methods for constraining Asian sources of CO using satellite (MOPITT) measurements of CO columns. *J. Geophys. Res.* **114**, D04305.
- Larssen, S., Sluyter, R. and Helmis, C. 1999. *Criteria for EUROAIRNET: The EEA Air Quality Monitoring and Information Network*. Technical Report, European Environment Agency.
- Le Dimet, F.-X. and Talagrand, O. 1986. Variational algorithms for analysis and assimilation of meteorological observations: theoretical aspects. *Tellus A* **38**, 97–110.
- Ménard, R., Cohn, S. E., Chang, L.-P. and Lyster, P. M. 2000. Assimilation of stratospheric chemical tracer observations using a Kalman filter. Part I: formulation. *Mon. Weather Rev.* **128**, 2654–2671.
- Mulholland, M. and Seinfeld, J. 1995. Inverse air pollution modelling of urban-scale carbon monoxide emissions. *Atmos. Environ.* **4**, 497–516.
- Müller, J.-F. and Stavrakou, T. 2005. Inversion of CO and NO_x emissions using the adjoint of the Image model. *Atmos. Chem. Phys.* **5**, 1157–1186.
- Nappo, C. J., Caneill, J. Y., Furman, R. W., Gifford, F. A., Kaimal, J. C. and co-authors. 1982. Workshop on the representativeness of meteorological observations. *B. Am. Meteorol. Soc.*, **63**, 761–764.
- Pétron, G., Granier, C., Khattatov, B., Lamarque, J.-F., Yudin, V. and co-authors. 2002. Inverse modeling of carbon monoxide surface emissions using Climate Monitoring and Diagnostics Laboratory network observations. *J. Geophys. Res.* **107**, 4761.
- Pétron, G., Granier, C., Khattatov, B., Yudin, V., Lamarque, J.-F. and co-authors. 2004. Monthly CO surface sources inventory based on the 2000–2001 MOPITT satellite data. *Geophys. Res. Lett.* **31**, L21107.
- Quélo, D., Krysta, M., Bocquet, M., Isnard, O., Minier, Y. and co-authors. 2007. Validation of the Polyphemus platform on the ETEX, Chernobyl and Algeciras cases. *Atmos. Environ.* **41**, 5300–5315.
- Saïde, P., Bocquet, M., Osses, A. and Gallardo, L. 2011. Constraining surface emissions of air pollutants using inverse modeling: method intercomparison and a new two-step multi-scale approach. *Tellus B* **63**, 360–370.
- Stavrakou, T. and Müller, J.-F. 2006. Grid-based versus big region approach for inverting CO emissions using Measurement of Pollution in the Troposphere (MOPITT) data. *J. Geophys. Res.* **111**, D15304.
- Vestreng, V., Myhre, G., Fagerli, H., Reis, S. and Terrasón, L. 2007. Twenty-five years of continuous sulphur dioxide emission reduction in Europe. *Atmos. Chem. Phys.* **7**, 3663–3681.
- Wu, L., Bocquet, M., Lauvaux, T., Chevallier, F., Rayner, P. and Davis, K. 2011. Optimal representation of source-sink fluxes for mesoscale carbon dioxide inversion with synthetic data. *J. Geophys. Res.* **116**, D21304.
- Yumimotoa, K. and Uno, I. 2006. Adjoint inverse modeling of CO emissions over Eastern Asia using four-dimensional variational data assimilation. *Atmos. Environ.* **40**, 6836–6845.
- Zhang, Y., Bocquet, M., Mallet, V., Seigneur, C. and Baklanov, A. In press. Real-time air quality forecasting, part II: state of the science, current research needs, and future prospects. *Atmos. Environ.* DOI:10.1016/j.atmosenv.2012.02.041.