

Bridging classical data assimilation and optimal transport: The 3D-Var case

Marc Bocquet¹, Pierre J. Vanderbecken¹, Alban Farchi¹, Joffrey Dumont Le Brazidec¹, and Yelva Roustan¹

¹CEREA, École des Ponts and EDF R&D, Île-de-France, France

Abstract. Because optimal transport acts as displacement interpolation in physical space rather than as interpolation in value space, it can ~~potentially avoid double penalty errors~~ avoid double-penalty errors generated by mislocations of geophysical fields. As such it provides a very attractive metric for non-negative ~~physical~~, sharp fields comparison — the Wasserstein distance — which could further be used in data assimilation for the geosciences. The algorithmic and numerical implementations of such distance are however not straightforward. Moreover, its theoretical formulation within typical data assimilation problems face conceptual challenges, resulting in scarce contributions on the topic in the literature.

We formulate the problem in a way that offers a unified view on both classical data assimilation and optimal transport. The resulting *OTDA* framework accounts for both the classical source of prior errors, background and observation, together with a Wasserstein barycentre in between states that stand for these background and observation. We show that the hybrid OTDA analysis can be decomposed as a simpler OTDA problem involving a single Wasserstein distance, followed by a Wasserstein barycentre problem which ignores the prior errors and can be seen as a *McCann interpolant*. We also propose a less enlightening but straightforward solution to the full OTDA problem, which includes the derivation of its analysis error covariance matrix. Thanks to these theoretical developments, we are able to extend the classical ~~3D-Var~~3D-Var/BLUE paradigm at the core of most classical data assimilation schemes. The resulting formalism is very flexible and can account for sparse, noisy observations and non-Gaussian error statistics. It is illustrated by simple one- and two-dimensional examples that show the richness of the new types of analysis offered by this unification.

1 Introduction

1.1 ~~Weakness of classical data~~ Data assimilation and the double-penalty issue

Geophysical data assimilation is a set of methods and algorithms at the intersection of Earth sciences, mathematics, and computer science, designed to enhance our understanding and predictive capabilities of the complex systems that govern our planet (Carrassi et al., 2018). These systems encompass the atmosphere, ocean, atmospheric chemistry and biogeochemistry, land surfaces, glaciology, hydrology, etc., and as a whole the climate system. Data assimilation is meant to optimally combine all sources of quantitative information, typically past and present observations, and numerical and statistical models of the system under consideration. Data assimilation (DA) is critical in forecasting chaotic geofluids by resetting the initial conditions

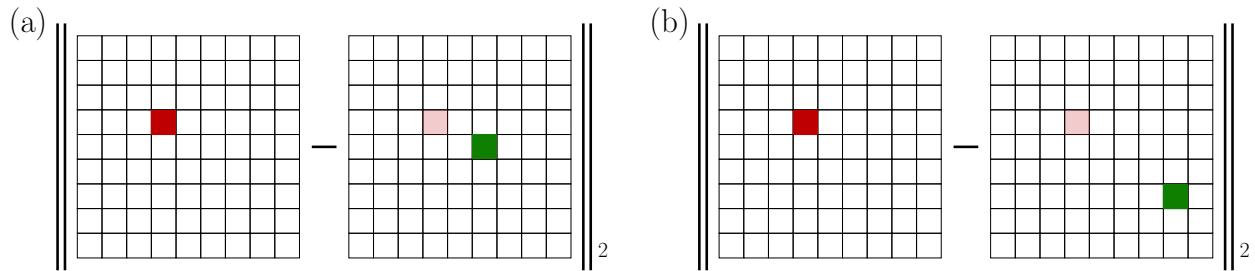


Figure 1. These two panels schematise the computation of the RMSE of two analysis increments. Those increments are the difference between the truth (left mesh within both Euclidean norm), concentrated here in the red grid-cell and the analysis located in the green grid-cells (right mesh within both Euclidean norm). The increment on panel (a) is the outcome of a better analysis spatially closer to the truth, as compared to the one in panel (b), and yet both increments yield the same RMSE. This verification metrics is hence impacted by the double-penalty error and does not help in discriminating location errors.

25 of the flow, estimating physical and statistical parameters of the models, and providing a quantitative re-analysis of the past history of the climate system over decades. Because classical DA is applied to complex and high-dimensional dynamics, the DA algorithms often result from a compromise between the sophistication of the employed mathematical techniques and their numerical scalability and efficiency (Kalnay, 2003; Asch et al., 2016; Evensen et al., 2022). For instance, it is well-known that most DA methods are built around or from an update step – the analysis – where observations and background states are
 30 combined, an operation which often relies on Gaussian statistical assumptions.

Here we would like to focus on ~~two other important weaknesses of classical DA. Firstly, one important issue that impacts classical DA, known as~~ the double-penalty error in ~~geosciences, which the geosciences. The double-penalty issue~~ refers to the over-penalisation of errors in both the model and observational data (e.g. Amodei and Stein, 2009) ~~and~~ compromises the balance required for effective DA. ~~It often stems from mislocation of fields which is caused by model error, in either the~~
 35 ~~forecasting or observation operator.~~ A typical example is given by a slight mislocation of a plume of pollutant resulting in high predicted concentration values at positions where no pollutant is observed while the model misses the observed peaks of concentration (Farchi et al., 2016). This mismatch is heavily penalised because of the use, over the same discretised space, of the root-mean square error (RMSE) utilised for a point-by-point comparison. Figure 1 shows an exemplar for such double-penalty error resulting in the inability to properly evaluate a model and learn from an analysis increment. This double-penalty
 40 error, a very common contribution to the *representation error* (Janjić et al., 2018), is ubiquitous in the geosciences: in numerical weather prediction and in particular for water vapour, in atmospheric chemistry ~~in and air quality, in biogeochemistry and in~~ eddy resolving ocean forecasting, etc. ~~This especially applies to sharp fields while it may be of less relevance for smoother, larger scale fields such as temperature.~~

~~A second weakness is the requirement in classical DA for prior fields, which could be the first-guess or background state and the observations, with substantial overlap in both space and time. For instance, failing to do so is a known pathway for the degeneracy of particle filters in high-dimension, also known as the curse of dimensionality (e.g. Farehi and Bocquet, 2018, and references th~~

~~But, even with Gaussian-based methods less prone to the curse of dimensionality, this can be seen as a drawback. Indeed, the update of classical DA would essentially interpolate in~~ It has been recognised that while the weighted Euclidean (Mahalanobis) distance can handle amplitude and smoothness mismatch, it cannot cope with mislocation error and hence account for the full *distortion* between mismatched fields (Hoffman et al., 1995). Hence, even though tuning covariances of Gaussian error distributions as in classical DA, such as increasing the correlation length, might help mitigate the double-penalty error, it is insufficient. In Fig. 1, one might replace the Euclidean norm by a weighted Euclidean one with a large correlation length. This would yield similar norm values for both cases. Unfortunately, it is not difficult to show that in this limit this (almost singular) norm can only distinguish between the spatial mean of both fields; it became blunt with no discriminating power. Towards DA, Fig. 1 in Feyeux et al. (2018) also illustrates why the Euclidean distance cannot properly cope with mislocation error. Note that, should Feyeux et al. (2018) have used a weighted Euclidean distance instead, with the same covariance matrix for the two contributions of the cost function, ~~the space of the values of the fields, yielding an analysis still confined within the support of the background state and that of the observation. This can be seen as a severe flaw if the mismatch in the observations and background state is due to an error in the location of the fields, or more generally when these fields are misspecified.~~ Figure ?? illustrates of couple of classical DA experiments with a beneficial analysis and a useless analysis resulting analysis state would have been the same and, in particular, independent from the covariance matrix. A similar but two-dimensional illustration is given by Fig. 3 in Vanderbecken et al. (2023).

~~These two panels illustrate a reasonably beneficial and a probably useless classical DA update. It is assumed that the observations (red dots) and first guess (dashed blue curve), which represent one-dimensional puffs of pollutant, are subject to a location mismatch. With a significant overlap of these fields, panel (a) displays a consistent analysis. By contrast, in panel (b), the analysis state (green curve) fails to propose a state with a significant presence in between the observations and first guess, just as expected since classical DA is based on an interpolation in field values, but neither in space, nor in time.~~

1.2 Nonlocal verification metrics

~~In both cases, the~~ The issue can be ~~ascribed~~ attributed to the use of a *local* verification metric, meaning that it compares, through the RMSE, values at the same site, of the same grid-cell. Thus, this issue goes beyond DA, and pertains to the use of local metrics.

To avoid being impacted by the double-penalty issue stemming from the use of local verification metrics, smarter *nonlocal* or *multiscale* metrics have been proposed. A typical metric of this kind consists in the combination of a displacement map followed by the use of classical norm such as the RMSE (Hoffman et al., 1995; Keil and Craig, 2009). In this vein, effective verification metrics can be based on *optical flow-based warping*, or on *deformed meshes*, prior to using classical norms (Gilleland et al., 2010a, b). These metrics can also be designed as *scale-dependent* and possibly *multiscale*, based on an empirical separation of scales, such as with *fuzzy* metrics (Ebert, 2008; Amodei and Stein, 2009), or e.g., *wavelets* (Briggs and Levine, 1997). They can be designed to grasp and quantify objects and features, such as lows and highs (Davis et al., 2006a, b; Lack et al., 2010). Metrics with similar capabilities but not necessary based on a displacement concept, have been introduced in computer vision

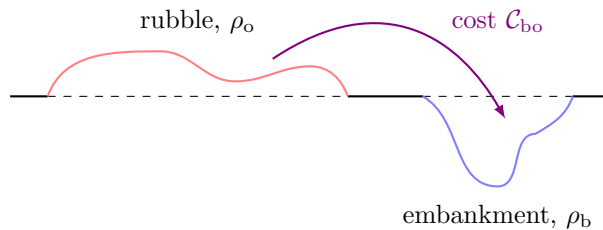


Figure 2. Illustration of the earth mover problem introduced by Monge in 1781 (see bulk of paper).

80 such as the *structural similarity index* (Zhou et al., 2004), or in the verification of precipitations (Skok, 2023; Necker et al., 2023).

One of the most elegant approach is based on the theory of *optimal transport* (OT), and the associated *Wasserstein distance*, which sits on solid mathematical foundations and significant developments, which are the main reasons why we will focus on OT in the following. Examples of application of OT to the verification of tracer and greenhouse gases models are given in
85 Farchi et al. (2016); Vanderbecken et al. (2023).

1.3 Optimal transport and the Wasserstein distance

Before mentioning applications of the Wasserstein distance in the field of geoscience, let us first give a very brief introduction to the concept and mathematical formulation of OT.

The OT concept stemmed from an engineering but rather universal problem. Gaspard Monge (Monge, 1781) considered the
90 *earth mover* problem where the goal is to efficiently move rubble to an embankment of about the same volume (see Fig. 2). Each displacement of a bit of earth has a known cost, so that the goal is to find the cheapest deterministic map that completely moves the rubble to the embankment. In mathematical terms, the goal is to find the map of minimal cost that transports the origin measure ρ_o to the target measure ρ_b , where *measure* here means that both of them are non-negative, and are integrable of integral 1. Note that the value 1 is arbitrary here and can be changed to $m > 0$, provided this is the mass of both ρ_o and ρ_b .
95 The cost is defined by a non-negative function \mathcal{C}_{bo} of two variables (one for the origin space and the other for the target space). Let us assume a quadratic cost, defined for any couple of points (x, y) of a geometric domain Ω :

$$\mathcal{C}_{bo}(x, y) = \|x - y\|_2^2, \quad (1)$$

where $\|\cdot\|_2$ is the Euclidean norm. Let us define the set of all admissible differentiable maps T that transport ρ_o to ρ_b :

$$\mathcal{U}_{bo} = \{T : \Omega \mapsto \Omega, \quad \rho_o = |\partial_x T| \rho_b \circ T\}, \quad (2)$$

100 where $|\partial_x T|$ is the absolute value of the determinant of the Jacobian of T , a factor which accounts for the deformation of the measure by the globally mass-conserving T . The *square* of the Wasserstein distance $\mathcal{W}_{\mathcal{C}_{bo}}$ is then defined by

$$\mathcal{W}_{\mathcal{C}_{bo}}^2(\rho_o, \rho_b) = \min_{T \in \mathcal{U}_{bo}} \int_{\Omega} \mathcal{C}_{bo}(x, T(x)) \rho_o(x) dx, \quad (3)$$

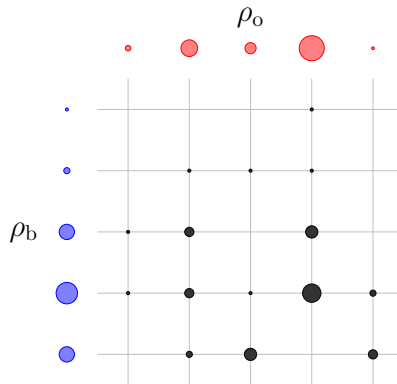


Figure 3. A representation of a discrete transport plan between two discrete origin (blue) and target (red) measures. The black dots represent the value of the transference plan. The radius of the dots are proportional to the values of these measures. This transference plan is checked to be admissible but is not necessarily optimal.

whose purpose is to minimise the total transport cost between ρ_o and ρ_b . It can be shown that $\mathcal{W}_{\mathcal{C}_{b_o}}$ is indeed a proper mathematical distance. The mathematical formulation is deceptively simple since it is elegant, compact and easy to grasp, but its theoretical and numerical solutions are far from trivial.

A breakthrough was made in the 20th century by Leonid Kantorovich who promoted the Monge problem to a probabilistic formulation. In his point of view, a bit of earth can be split and moved to many sites of the target measure support. The deterministic map T is hence replaced with a probabilistic measure π defined over $\Omega \times \Omega$. Such a π is called hereafter a *transference plan*. An admissible transference plan is integrable and have ρ_o and ρ_b as one-variable marginals; hence the definition of the admissible set:

$$\mathcal{V}_{b_o} = \left\{ \pi : \Omega \times \Omega \mapsto \mathbb{R}_+, \quad \rho_o(x) = \int_{\Omega} \pi(x, y) dy, \quad \rho_b(y) = \int_{\Omega} \pi(x, y) dx \right\}. \quad (4)$$

As opposed to the deterministic Monge maps, the transference plans offer a symmetrical view on the origin and target space and their measures. An illustration of a discrete transference plan is given by Fig. 3. In this view the squared Wasserstein distance can be reformulated as

$$\mathcal{W}_{\mathcal{C}_{b_o}}^2(\rho_o, \rho_b) = \min_{\pi \in \mathcal{V}_{b_o}} \int_{\Omega \times \Omega} \mathcal{C}_{b_o}(x, y) \pi(x, y) dx dy. \quad (5)$$

Equations (3,5) are the consecrated continuous formulations of OT. In the rest of the paper, we will deal instead with *discrete* related formulations, more tangible, and amenable to algorithmic and numerical implementations.

The field has attracted a lot of attention from pure and applied mathematicians, and computer scientists. A complete introduction to the topic by its experts can be found in the stimulating text books by Vilani (2003, 2009); Peyré and Cuturi (2019). Peyré and Cuturi (2019) nicely provide concrete examples, numerical methods and a broad coverage of the topic from the perspective of applied mathematicians and computer scientists. Hence, it will be referred to quite often in the rest of the paper.

1.4 Nonlocal, multiscale metrics and data assimilation

Let us now go back to DA and narrow our focus to the use of advanced metrics in DA. Accounting for displacement error in DA and hence relying on nonlocal verification metrics has been advocated by [Ravela et al. \(2007\)](#); [Plu \(2013\)](#); [Hoffman and Grassotti \(1996\)](#); [Ravela et al. \(2007\)](#). Metrics built on a multiscale analysis of the fields to achieve a similar goal have been proposed by Ying (2019); Ying et al. (2023).

Wasserstein distance and closely related formulations, have been advocated in the flow formulation of the analysis (DA update) to seamlessly transport the prior to the posterior (El Moselhy and Marzouk, 2012; Oliver, 2014; Marzouk et al., 2017; Farchi and Bocquet, 2018; Tamang et al., 2020). It can for instance be used to adjust the posterior discrete probability density functions (pdf) in the particle filter. It has similarly been used to assist ensemble DA (Tamang et al., 2021, 2022). Finally, it has also very recently been used to compare forecast ensembles for sub-seasonal prediction (Le Coz et al., 2023; Lledó et al., 2023), or precipitation ([?](#)) ([Duc and Sawada, 2024](#)).

In the context of this paper, it is *critical* to be aware that the use of OT in practical DA focused so far on applying OT to the pdf of a single variable. Quite often, OT is applied to the pdf of a single random variable because

- OT in one dimension (the space of the values taken by this random variable) together with the quadratic cost has a very simple solution that only relies on the cumulative distribution functions of the origin and target measures (see e.g., Remark 2.30 in Peyré and Cuturi, 2019).
- increasing the number of random variables is subject to the curse of dimensionality, necessitating an exponential increase in computational resources, [when increasing the resolution of the discretised fields](#).

This is very different from our context and objective where the objects dealt with by OT are [\(non-negative\) physical](#) field states, not the pdf of one of their variables. In particular, although the computations are very demanding when the physical space that supports the fields scales to dimension 2 and 3, our problem is not subject to the curse of dimensionality.

1.5 [Feyeux et al. proposal](#)

The present paper stands more in the wake of the seminal [proposal of proposals of Ning et al. \(2014\) and Feyeux \(2016\)](#); [Feyeux et al. \(2018\)](#). Their idea is to replace the local metrics of classical variational DA, typically the square of the Euclidean distance (hence related to the L_2 norm) by the squared Wasserstein distance. This is intuitively what we are after in order to cope with [the two weaknesses mislocation errors](#) mentioned in Sect. 1.1 in the context of DA. This should redefine the nature of the DA update step. Let us formalise this idea ([Feyeux, 2016](#)).

We will seize this opportunity to introduce some of our notation, in the context of discrete DA which is a widely adopted standpoint in the geosciences. Let us focus on a classical DA ~~3D-Var~~ [3D-Var](#) cost function (Daley, 1991):

$$G_{cl}(\mathbf{x}^a) = \|\mathbf{y}^b - \mathbf{x}^a\|_2^2 + \|\mathbf{y}^o - \mathbf{H}\mathbf{x}^a\|_2^2, \quad (6)$$

where $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^N x_i^2}$ is the Euclidean norm, $\mathbf{y}^b \in \mathbb{R}^{N_b}$ is the first guess, $\mathbf{y}^o \in \mathbb{R}^{N_o}$ is the vector of observations, and \mathbf{H} is the observation operator¹. $\mathbf{x}^a \in \mathbb{R}^{N_a}$ is the dummy variable of this optimisation problem whose optimal value corresponds to the DA state analysis. Now, the substitution of the Euclidean norm yields the new 3D-Var-3D-Var cost function:

$$155 \quad G_w(\mathbf{x}^a) = \mathcal{W}_2^2(\mathbf{y}^b, \mathbf{x}^a) + \mathcal{W}_2^2(\mathbf{y}^o, \mathbf{H}\mathbf{x}^a), \quad (7)$$

where \mathcal{W}_2 is some discretisation of the Wasserstein distance based on the cost defined by the square of the Euclidean distance. Note that this 3D-Var-3D-Var requires balancing two instances of a Wasserstein-based metric. The analysis state is known as a *Wasserstein barycentre*, abridged W-barycentre in the following.

Feyeux (2016); Feyeux et al. (2018) explored the optimisation aspects of this DA problem. However, Feyeux (2016) ultimately pointed to a possible inconsistency in the definition of the DA problem formulated in Eq. (7), where the system is only partially observed (non-trivial \mathbf{H}). In the case where the system is fully observed, typically when \mathbf{H} is the identity operator, the outcome of the optimisation problem, i.e. the analysis, matches our expectations. However, when the system is partially observed, inconsistencies are observed. Let us see why.

Panel (a) of Fig. 4 considers the DA problem based on Eq. (7), assuming that only half of the domain is observed. We have solved the corresponding mathematical and numerical problem as raised by Feyeux (2016) and displayed its solution. However, most-of-one observes that the mass of the solution concentrates on the observed subdomain, and neglects the rest of the domain where the prior mainly concentrates, an outcome suspected by Feyeux (2016). Instead, we would have anticipated intuitively preferred a solution close to the one offered by panel (b) of Fig. 4, whose formulation and numerical solution differ and follow the new approach developed in the present paper (how we obtained this solution will be described in Sect. 2).

The main caveat of Eq. (7) comes from the implicit-assumption-occultation of part of the domain (the kernel of \mathbf{H} is non-trivial), together with the requirement that OT is *balanced*, i.e. the origin and target densities must have the same mass. This mass balance applies to both OTs-OT terms in Eq. (7), between \mathbf{y}^b and \mathbf{x}^a and between \mathbf{y}^o and $\mathbf{H}\mathbf{x}^a$:

$$175 \quad m(\mathbf{x}^a) = m(\mathbf{y}^b), \quad m(\mathbf{H}\mathbf{x}^a) = m(\mathbf{y}^o), \quad (8)$$

where the mass of the-a vector $\mathbf{x} \in \mathbb{R}^N$ is defined by

$$175 \quad m(\mathbf{x}) = \mathbf{1}^\top \mathbf{x} = \sum_{i=1}^N x_i, \quad (9)$$

with $\mathbf{1} \in \mathbb{R}^N$ hereafter defined as the vector of entries 1.

Now, if we further assume for simplicity that \mathbf{y}^b and \mathbf{y}^o have the same mass (which is the case in Fig. 4), then

$$m(\mathbf{H}\mathbf{x}^a) = m(\mathbf{y}^o) = m(\mathbf{y}^b) = m(\mathbf{x}^a). \quad (10)$$

As a result, the-mass-equality-we obtain $m(\mathbf{H}\mathbf{x}^a) = m(\mathbf{x}^a)$ compels-to-find-part-which-is-an-undesired-prior-piece-of-information-as-to-where-to-find-the-mass of \mathbf{x}^a in $\mathbb{R}^{N_a} \setminus \ker(\mathbf{H})$, in-order-to-fully-account-for- $m(\mathbf{x}^a)$. Hence-. Simply put, unless the system

¹The notation \mathbf{y}^b and \mathbf{y}^o is at variance with the more familiar \mathbf{x}^b and \mathbf{y} notation of DA, respectively. Yet, this change will prove very useful in the following: it follows the idea that the full *information vector* is $\mathbf{y} = [(\mathbf{y}^b)^\top, (\mathbf{y}^o)^\top]^\top$ whose components may benefit from homogeneous notation (Talagrand, 1997).

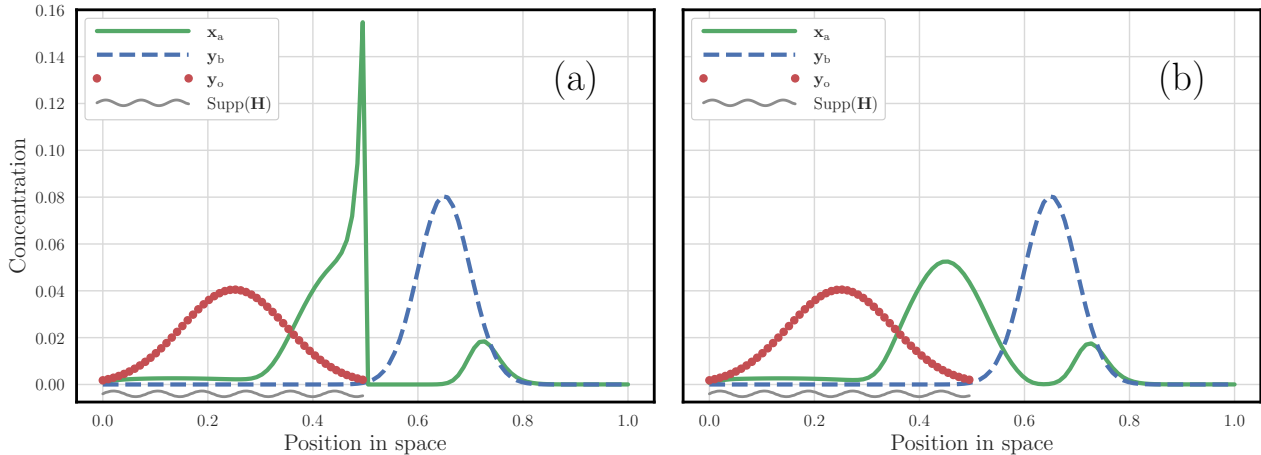


Figure 4. These panels illustrate the analysis of a ~~3D-Var~~3D-Var that relies on the Wasserstein distance rather than a local metric. The red dots represent the observations, while the dashed blue curve represents the background state. The observations are only focused on the left half of the domain. The solution of the optimisation problem Eq. (7), is displayed on panel (a) as a solid green curve. The solution of the optimisation problem we will propose in this paper is displayed on panel (b) as a solid green curve. The support of the observation is suggested with a wavy grey segment. These states are typically one-dimensional puff pollutant concentrations. They should not be confused with pdfs of a single random variable.

is fully observed, this approach amounts to *finding the truth under the streetlight*, ~~unless the system is fully observed~~. This is precisely what happens in panel (a) of Fig. 4 with the undesired concentration of the mass of x^a close to the edge of the observed subdomain.

To overcome this caveat and find a proper alternative to Eq. (7), we need ~~to resort to (i) to renounce comparing the fields in~~
 185 observation space (in the observation discrepancy term of the cost function), and (ii) introduce unbalanced OT, i.e. we need to be able to accommodate states of distinct masses. In the computer science context of pure OT, ~~such possibility the latter~~ has been discussed by Chizat et al. (2018). But our solution differs formally and will be DA-centric.

1.6 Objective and outline

The objective of this paper is to lift the objection of Feyeux (2016), and propose a DA framework based on the Wasserstein
 190 distance, and hence to offer a consistent way to bridge OT and classical DA. The new formalism will be referred to as *hybrid OTDA* for *hybrid optimal transport data assimilation* in the rest of this paper, and often for the sake of brevity *OTDA*. We will ~~mainly~~ focus on the definition of a ~~3D-Var~~3D-Var DA problem and how to obtain its analysis state and the associated analysis error covariance matrix.

At least within the perimeter of this paper, some restrictions apply ~~compared to traditional DA~~. Firstly, the physical fields
 195 considered in the DA problem are *non-negative* (concentration of tracer, pollutants, water vapour, hydrometeors, chemical and

biogeochemical species, etc.). However, as opposed to Feyeux (2016), the methods of this paper do not require the (possibly noisy) background state \mathbf{y}^b and observation \mathbf{y}^o to be non-negative. We stress once again that the states of our DA problem are physical fields onto which OT is applied and are not meant to be pdf of a random variable. Secondly, the observation operator \mathbf{H} is assumed to be *linear*. This is only meant for convenience and to obtain a rigorous correspondence between the primal and dual cost functions of the ~~3D-Var~~. Finally, we stress once again that the states of our DA problem are physical fields onto which OT is applied and are not meant to be pdf of a random variable ~~3D-Var~~. Making this assumption is very common in geophysical DA: \mathbf{H} can indeed be seen as the tangent linear of a nonlinear observation operator within the inner loop of a 3D-Var or a 4D-Var (see for instance Courtier, 1997).

The outline of the paper follows. After the present introduction, Sect. 2 discloses our main idea, and discusses two mathematical paths to solve the underlying optimisation problem, a first one which is enlightening but not necessarily practical and an alternative which is direct and robust but hides some of the concepts behind it. Section 3 provides one- and two- dimensional illustrations of a ~~3D-Var~~ 3D-Var analysis based on the new hybrid OTDA formalism. These illustrations will show the possibilities and flexibility of the new framework. Importantly, This section will also depict classical DA as limit case of the formalism. In Sect. 4, the second-order analysis, i.e. the uncertainty quantification of the OTDA ~~3D-Var~~ 3D-Var, is derived, discussed, and illustrated. Conclusions and perspectives are given in Sect. 5.

2 The main proposal

2.1 Notation and conventions

Non-negative vectors \mathbf{x} of size N are called *discrete measures*; they lie in the orthant $\mathcal{O}_N^+ \triangleq \mathbb{R}_+^N$. Although most mathematical OT theories work on *normalised* discrete measures, yielding *probability vectors*, this assumption won't be needed in this paper. The open subset of \mathcal{O}_N^+ of all the positive discrete measures will be denoted $\mathcal{O}_{N,+}^{+,*} \triangleq \mathbb{R}_{+,*}^N$.

We will distinguish the *observations* $\mathbf{y}^b \in \mathbb{R}^{\mathfrak{N}_b}$, $\mathbf{y}^o \in \mathbb{R}^{\mathfrak{N}_o}$ from the observable states $\mathbf{x}^b \in \mathcal{O}_{N_b}^+$ and $\mathbf{x}^o \in \mathcal{O}_{N_o}^+$. \mathbf{y}^b , which corresponds to the first guess of conventional DA, and \mathbf{y}^o , which corresponds to the traditional observation vector, are known before solving the ~~3D-Var~~ 3D-Var problem. These vectors embody all the information processed in the analysis. By contrast, the observables \mathbf{x}^b and \mathbf{x}^o , which are related to \mathbf{y}^b and \mathbf{y}^o , respectively, through an observation operator (the identity for \mathbf{y}^b and \mathbf{x}^b), are not known a priori. They will be estimated together with the analysis state $\mathbf{x}^a \in \mathcal{O}_{N_a}^+$. Note that these vectors may well lie in distinct vector spaces of different dimensions, hence the introduction of as many dimensions $N_b, N_o, \mathfrak{N}_b, \mathfrak{N}_o$. x_i^* can be seen as the value taken by \mathbf{x}^* at site \mathbf{r}_i^* , for $\star = b, o, a$ and $i \in \llbracket 1, N_\star \rrbracket$. Mind that the distinction between \mathbf{y}^b and \mathbf{x}^b , and the introduction of \mathbf{x}^o is a novelty of OTDA compared to classical DA.

Like in classical DA, the vectors \mathbf{y}^b and \mathbf{y}^o are subject to (prior) errors whose statistics are specified by the likelihoods $p(\mathbf{y}^b|\mathbf{x}^b)$ and $p(\mathbf{y}^o|\mathbf{x}^o)$, respectively. Up to constants that do not depend on $\mathbf{x}^b, \mathbf{x}^o, \mathbf{y}^b, \mathbf{y}^o$, we assume the existence of ζ_b and ζ_o such that

$$\ln p(\mathbf{y}^b|\mathbf{x}^b) \triangleq -\zeta_b(\mathbf{y}^b - \mathbf{x}^b) + \text{cst}, \quad \ln p(\mathbf{y}^o|\mathbf{x}^o) \triangleq -\zeta_o(\mathbf{y}^o - \mathbf{H}\mathbf{x}^o) + \text{cst}, \quad (11)$$

so that various error statistics can be considered. These errors are hypothesised to be mutually independent. The observation operator $\mathbf{H} : \mathcal{O}_{N_o}^+ \mapsto \mathbb{R}^{\mathfrak{M}_o}$ used in the definition of ζ_o is assumed to be linear. This qualification is for convenience and could be lifted if necessary. It is further assumed that ζ_b and ζ_o are strictly convex functions. This is for instance the case if we choose Gaussian error statistics yielding

$$\zeta_b(\mathbf{e}_b) = \frac{1}{2} \|\mathbf{e}_b\|_{\mathbf{B}^{-1}}^2, \quad \zeta_o(\mathbf{e}_o) = \frac{1}{2} \|\mathbf{e}_o\|_{\mathbf{R}^{-1}}^2, \quad (12)$$

where $\|\mathbf{e}\|_{\mathbf{A}} = \sqrt{\mathbf{e}^\top \mathbf{A} \mathbf{e}}$. \mathbf{B} is the positive definite background error covariance matrix, and \mathbf{R} is the positive definite observation error covariance matrix. Finally, the $m(\star)$ operator will act in the following on not only vectors but, more generally, on any tensor and will return the sum of all of its entries.

2.2 Formalism of discrete optimal transport

To discretise and solve the continuous Kantorovich optimisation problem introduced in Sect. 1.3, we will need two elementary pieces of information about OT. These are not the only techniques we will leverage, but both represent cornerstones towards a numerical solution to our proposal, and hence they need a proper introduction.

2.2.1 The primal cost function

Let us consider two discrete measures $\mathbf{x}^b \in \mathcal{O}_{N_b}^+$ and $\mathbf{x}^o \in \mathcal{O}_{N_o}^+$ having the same mass

$$m \triangleq m(\mathbf{x}^b) = m(\mathbf{x}^o). \quad (13)$$

For convenience, $\mathcal{O}_{b,o}^+$ will be used as an alias for the set $\mathcal{O}_{N_b \times N_o}^+$. A cost matrix $\mathbf{C}_{b,o} \in \mathcal{O}_{b,o}^+$ is given. The optimisation problem will be formulated using discrete Kantorovich transference plans $\mathbf{P}^{b,o} \in \mathcal{O}_{b,o}^+$. The optimal discrete transference plan is given by the minimiser of the following optimisation problem:

$$W_{\mathbf{C}_{b,o}}(\mathbf{x}^b, \mathbf{x}^o) \triangleq \min_{\mathbf{P}^{b,o} \in \mathcal{U}_{b,o}(\mathbf{x}^b, \mathbf{x}^o)} \text{Tr}(\mathbf{C}_{b,o}^\top \mathbf{P}^{b,o}), \quad (14a)$$

where the trace sums up the costs attached to each path, and the set of admissible transference plans is defined by

$$\mathcal{U}_{b,o} \triangleq \left\{ \mathbf{P} \in \mathcal{O}_{b,o}^+ : \mathbf{P} \mathbf{1}_o = \mathbf{x}^b, \quad \mathbf{P}^\top \mathbf{1}_b = \mathbf{x}^o \right\}, \quad (14b)$$

which selects the discrete transference plans with the proper marginals. $W_{\mathbf{C}_{b,o}}$ could be viewed as a discrete equivalent to the square of the Wasserstein distance $\mathcal{W}_{\mathbf{C}_{b,o}}^2$ introduced in Eq. (5).

2.2.2 Entropic regularisation

Adding to the fact that the optimisation problem Eq. (14a) is constrained, it may neither be convex, nor exhibit a single minimum. ~~Hence, entropic~~ Entropic regularisation is used addresses these issues and is used here to lift the constraints and

to render the problem strictly convex. A comprehensive justification is given by Peyré and Cuturi (2019). ~~Note however, that~~
 255 More precisely, we will use a *Kullback-Leibler divergence* (KL) regularisation term to be inserted in Eq. (14a),

$$\text{Tr}(\mathbf{C}_{\text{bo}}^\top \mathbf{P}^{\text{bo}}) \rightarrow \text{Tr}(\mathbf{C}_{\text{bo}}^\top \mathbf{P}^{\text{bo}}) + \varepsilon \mathcal{K}(\mathbf{P}^{\text{bo}} | \boldsymbol{\nu}^{\text{bo}}), \quad (15)$$

which incorporates some prior transference plan $\boldsymbol{\nu}^{\text{bo}}$ and does not require $m(\mathbf{P}^{\text{bo}}) = 1$, whereas Peyré and Cuturi (2019) opted for a basic entropy term. The KL term (Boyd and Vandenberghe, 2004) is defined by

$$\mathcal{K}(\mathbf{p} | \mathbf{q}) \triangleq \sum_i p_i \ln \frac{p_i}{q_i} - p_i + q_i. \quad (16)$$

260 It can be checked that the Hessian of the regularised cost function Eq. (15) is a diagonal matrix of coefficients $\varepsilon / P_{ij}^{\text{bo}} \geq \varepsilon$ since $0 \leq P_{ij}^{\text{bo}} \leq 1$, making the problem ε -strongly convex. We choose, e.g., $\boldsymbol{\nu}^{\text{bo}} = \mathbf{x}^{\text{b}} (\mathbf{x}^{\text{o}})^\top / m$, and $\varepsilon > 0$ which is the regularisation scalar parameter. Note that this particular $\boldsymbol{\nu}^{\text{bo}}$ is an admissible transference plan, i.e. it belongs to \mathcal{U}_{bo} , and can be interpreted as a complete statistical decoupling of the transference plan with respect to the origin and target discrete measures. In the limit $\varepsilon \rightarrow 0^+$ of vanishing regularisation, the solution should not depend on the choice of $\boldsymbol{\nu}^{\text{bo}}$. However, the
 265 convergence to the solution at finite ε may depend on this choice. The primal cost function augmented with such an entropic regularisation is usually solved numerically using the iterative *Sinkhorn algorithm* (Sinkhorn, 1964). However, this is not the path followed in this paper, although we have used it as well.

Finally note that the technique to convexify such an optimisation problem with a KL term has been introduced in DA by Bocquet (2009); Bocquet et al. (2011) following principles of statistical physics.

270 2.3 From classical data assimilation to hybrid optimal transport data assimilation

Figure 5 is a schematic representation of the flow of information in a classical DA update, and in particular ~~3D-Var~~3D-Var, using the notation introduced above. In this case, the observables \mathbf{x}^{b} , \mathbf{x}^{o} , and the analysis state \mathbf{x}^{a} ~~coincide~~are the same by construction, hence \mathbf{x}^{b} and \mathbf{x}^{o} are not needed. This diagram, which could also be seen as a Bayesian network, corresponds to the cost function

$$275 L_{\text{cl}}(\mathbf{x}^{\text{a}}) = \zeta_{\text{b}}(\mathbf{y}^{\text{b}} - \mathbf{x}^{\text{a}}) + \zeta_{\text{o}}(\mathbf{y}^{\text{o}} - \mathbf{H}\mathbf{x}^{\text{a}}), \quad (17)$$

to be minimised over \mathbf{x}^{a} . Now let us make use of the observables \mathbf{x}^{b} and \mathbf{x}^{o} as new degrees of freedom but bind them by OTs to \mathbf{x}^{a} , using the cost matrices \mathbf{C}_{ba} and \mathbf{C}_{oa} , respectively. This yields the diagram in Fig. 6, which corresponds to the cost function

$$L_{\text{w}}(\mathbf{x}^{\text{a}}) = \min_{\mathbf{x}^{\text{b}} \in \mathcal{O}_{N_{\text{b}}}^+, \mathbf{x}^{\text{o}} \in \mathcal{O}_{N_{\text{o}}}^+} \{ \zeta_{\text{b}}(\mathbf{y}^{\text{b}} - \mathbf{x}^{\text{b}}) + \zeta_{\text{o}}(\mathbf{y}^{\text{o}} - \mathbf{H}\mathbf{x}^{\text{o}}) + W_{\mathbf{C}_{\text{ba}}}(\mathbf{x}^{\text{b}}, \mathbf{x}^{\text{a}}) + W_{\mathbf{C}_{\text{oa}}}(\mathbf{x}^{\text{o}}, \mathbf{x}^{\text{a}}) \}. \quad (18)$$

280 It must be minimised over \mathbf{x}^{a} , yielding an analysis state \mathbf{x}^{a} which can also be seen as the W-barycentre between \mathbf{x}^{b} and \mathbf{x}^{o} . Note that \mathbf{x}^{b} and \mathbf{x}^{o} are discrete measures of unknown mass. For the optimisation problem, they lie in $\mathcal{O}_{N_{\text{b}}}^+$ and $\mathcal{O}_{N_{\text{o}}}^+$, respectively.

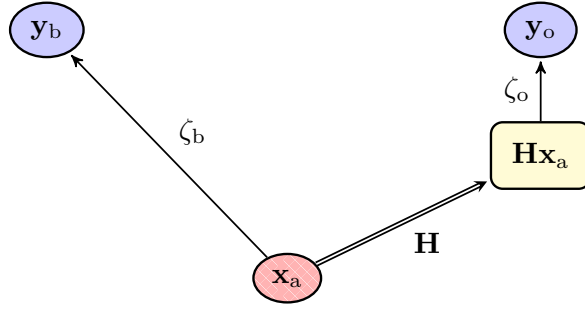


Figure 5. A diagrammatic representation of the classical $3D\text{-Var-}3D\text{-Var}$ update, with the observations \mathbf{y}^b (the first guess) and \mathbf{y}^o (the observation vector), the analysis state \mathbf{x}^a , and the observed analysis $\mathbf{H}\mathbf{x}^a$. A double line arrow represents a deterministic map, whereas a single line arrow represents a statistical binding between the origin and the target.

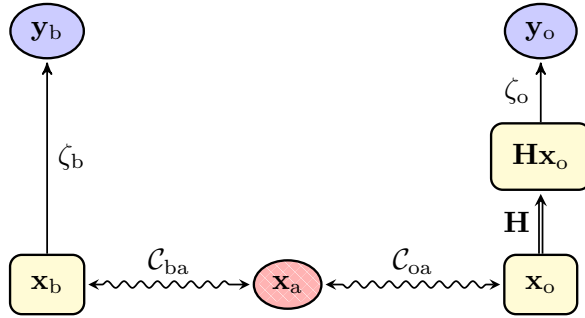


Figure 6. A diagrammatic representation of the hybrid OTDA $3D\text{-Var-}3D\text{-Var}$ update, with the observations \mathbf{y}^b (the first guess) and \mathbf{y}^o (the observation vector), the observables \mathbf{x}^b , \mathbf{x}^o , and \mathbf{x}^a which is the W-barycentre. A double line arrow represents a deterministic map, a single line arrow represents a statistical binding between the origin and the target, and a waving line represents the *weaker* bindings of \mathbf{x}^b and \mathbf{x}^o and \mathbf{x}^o and \mathbf{x}^a through OTs. This diagram can be seen as an unfolding of that of Fig. 5.

285 Moving from Eq. (17) to Eq. (18) following the principles and guidance of the introductory Sect. 1.4 is empirical, but no more than in Ning et al. (2014); Feyeux (2016). Showing its merits is the goal of the present paper. As opposed to Feyeux et al. (2018), it can deal with sparse and noisy observations, i.e. non-trivial \mathbf{H} . We will show that classical DA is embedded in this generalisation. Moreover, the merits of the new cost function will be a posteriori qualitatively supported by the outcome of the numerical experiments (to the expert's eyes), which improve over previous formalism's outcomes. We would like to point out that we have also developed a consistent probabilistic and Bayesian formalism fully supporting the introduction of Eq. (18). However we felt that the derivation is too long and technical for this paper, and would not be helpful in the exploration of the direct consequences of Eq. (18).

290

We call Eq. (18) a *high-level* primal cost function because the metrics $W_{C_{ba}}$ and $W_{C_{oa}}$ have not yet been replaced by their transference plan expression as opposed to, e.g., Eq. (14a). Passing to a lower level primal cost function would require to expand Eq. (18) using Eq. (14a) twice.

In the subsequent two subsections, we will investigate two pathways to solve the optimisation problem Eq. (18). The first path, Sect. 2.4, unveils some of the key concepts behind its solution, and partially disentangle the classical DA part from the W-barycentre part of the full analysis. This approach is enlightening but not necessarily practical. The second path is an alternative which is direct and robust but hides some of the fundamental principles underlying the solution. The busy reader could skip directly to the latter, i.e. Sect. 2.5.

2.4 Decomposition of the optimisation problem and effective cost metric

In this section, key ideas behind the minimisation of Eq. (18) are sketched and discussed. The level of mathematical rigour of this section is that of casual methodological DA in the geoscience literature, ~~not at the level of applied mathematics~~. However, we stress that all the algorithms discussed here have been tested numerically successfully on various configurations. The solution of Eq. (18) presented in this section is not necessarily robust, but it is enlightening and hence worth discussing.

Repeated contravariant indices – meaning the same tensor index present as upper and lower index – in tensor expressions will be understood as summed over, following Einstein’s convention.

2.4.1 Dual formulation of the full-primal problem

One way, though not the only one, to write the explicit primal problem associated to Eq. (18) is through the use of a *gluing* transference plan $\mathbf{P}^{boa} \in \mathcal{O}_{b,o,a}^+$ where $\mathcal{O}_{b,o,a}^+ = \mathbb{R}_+^{N_b N_o N_a}$ (see p.11-12 of Vilani, 2009), a 3–tensor whose marginals are \mathbf{x}^b , \mathbf{x}^o , and \mathbf{x}^a and that glues the transference plans \mathbf{P}^{ba} between \mathbf{x}^b and \mathbf{x}^a and \mathbf{P}^{oa} between \mathbf{x}^o and \mathbf{x}^a :

$$\mathcal{L} = \min_{\mathbf{x}^a \in \mathcal{O}_a^+} L_w(\mathbf{x}^a), \quad (19a)$$

$$= \min_{\mathbf{x}^b \in \mathcal{O}_b^+ \mathbf{x}^o \in \mathcal{O}_o^+ \mathbf{x}^a \in \mathcal{O}_a^+} \left[\zeta_b(\mathbf{y}^b - \mathbf{x}^b) + \zeta_o(\mathbf{y}^o - \mathbf{H}\mathbf{x}^o) + \min_{\mathbf{P} \in \mathcal{U}_{boa}} \{P_{ijk} C_{ba}^{ik} + P_{ijk} C_{oa}^{jk}\} \right]. \quad (19b)$$

where the admissible set of (glued) transference plans, the set of all 3–tensors of non negative entries whose marginals are \mathbf{x}^b , \mathbf{x}^o and \mathbf{x}^a , is defined by

$$\mathcal{U}_{boa} \triangleq \left\{ \mathbf{P} \in \mathcal{O}_{b,o,a}^+ : \forall i, P_{ijk} 1_o^j 1_a^k = x_i^b, \quad \forall j, P_{ijk} 1_b^i 1_a^k = x_j^o, \quad \forall k, P_{ijk} 1_b^i 1_o^j = x_k^a \right\}. \quad (19c)$$

Because of the hardly scalable dimensionality of the primal problem, based on either a 3–tensor, or a couple of 2–tensors, we wish to derive a dual problem equivalent to the primal one, using Lagrange multipliers to lift the constraints with, as will be checked later, a significantly smaller dimensionality.

This leads to ~~the following series of transformation~~ a series of transformations of the problem \mathcal{L} , from a Lagrangian to a dual cost function: ~~-, which is reported in Appendix A for the mathematics-inclined reader. The outcome is a dual problem which~~

$$\mathcal{L}^* = \max_{(\mathbf{f}_b, \mathbf{f}_o) \in \mathcal{U}_{b_o}^*(\mathbf{C}_{ba}, \mathbf{C}_{oa}, \mathbf{H})} \{ \mathbf{f}_b^\top \mathbf{y}^b + \mathbf{f}_o^\top \mathbf{y}^o - \zeta_b^*(\mathbf{f}_b) - \zeta_o^*(\mathbf{f}_o) \}, \quad (20a)$$

where the $*$ symbol refers to *dual* and where the polyhedron $\mathcal{U}_{b_o}^*(\mathbf{C}_{ba}, \mathbf{C}_{oa}, \mathbf{H})$ is defined by

$$\mathcal{U}_{b_o}^*(\mathbf{C}_{ba}, \mathbf{C}_{oa}, \mathbf{H}) \triangleq \left\{ \mathbf{f}_b \in \mathbb{R}^{\mathfrak{N}_b}, \mathbf{f}_o \in \mathbb{R}^{\mathfrak{N}_o} : \forall i, j, k, \quad f_b^i + f_o^l H_l^j \leq C_{ba}^{ik} + C_{oa}^{jk} \right\}. \quad (20b)$$

In Eq. (20), the maps ζ_b^* and ζ_o^* are the Legendre-Fenchel transforms of the maps ζ_b and ζ_o , respectively. Let us recall that the Legendre-Fenchel transform $\mathbf{f} \mapsto \zeta^*(\mathbf{f})$ of the map $\mathbf{e} \mapsto \zeta(\mathbf{e})$ is defined by $\zeta^*(\mathbf{f}) = \sup_{\mathbf{e}} \{ \mathbf{f}^\top \mathbf{e} - \zeta(\mathbf{e}) \}$. For instance, in the case of Gaussian error statistics as in Eq. (12), these transforms are given by

$$\zeta_b^*(\mathbf{f}_b) = \frac{1}{2} \|\mathbf{f}_b\|_{\mathbf{B}}^2, \quad \zeta_o^*(\mathbf{f}_o) = \frac{1}{2} \|\mathbf{f}_o\|_{\mathbf{R}}^2. \quad (21)$$

~~From Eq. to Eq., taking the minimum over the observables \mathbf{x}^b , \mathbf{x}^o , and \mathbf{x}^a implies to enforce $\mathbf{h}_b = \mathbf{f}_b$, $\mathbf{h}_o = \mathbf{H}^\top \mathbf{f}_o$, and $\mathbf{f}_a = \mathbf{0}$. Hence, we finally obtain the dual problem which only depends on the Lagrange multipliers: The inequality constraints of the polyhedron $\mathcal{U}_{b_o}^*$ stem from the positivity constraint $P_{ijk} \geq 0$ in Eq.. Very importantly, we have the coincidence of the minimum of the primal problem with the maximum of the dual problem $\mathcal{L} = \mathcal{L}^*$, a property called *strong duality* (see Sect. 5.2 in Boyd and Vandenberghe, 2004). Strong duality can for instance be achieved if both the primal and dual cost functions are convex, which is the case here. This allows~~

2.4.2 [Decomposition of the dual problem](#)

[These transformations allow](#) us to trade the primal for the dual problem. Since for each pair $\mathbf{f}_b, \mathbf{f}_o$ in $\mathcal{U}_{b_o}^*$, there are N_a constraints indexed by $k \in \llbracket 1, N_a \rrbracket$, and that the tightest of these constraints can account for the others, the problem Eq. (20) should be equivalent to

$$\mathcal{L}^* = \max_{(\mathbf{f}_b, \mathbf{f}_o) \in \mathcal{U}_{b_o}^*(\mathbf{C}_{bo}, \mathbf{H})} \{ \mathbf{f}_b^\top \mathbf{y}^b + \mathbf{f}_o^\top \mathbf{y}^o - \zeta_b^*(\mathbf{f}_b) - \zeta_o^*(\mathbf{f}_o) \}, \quad (22a)$$

where the polyhedron $\mathcal{U}_{b_o}^*(\mathbf{C}_{bo}, \mathbf{H})$ is defined by

$$\mathcal{U}_{b_o}^*(\mathbf{C}_{bo}, \mathbf{H}) \triangleq \left\{ \mathbf{f}_b \in \mathbb{R}^{\mathfrak{N}_b}, \mathbf{f}_o \in \mathbb{R}^{\mathfrak{N}_o} : \forall i, j, \quad f_b^i + f_o^l H_l^j \leq C_{bo}^{ij} \right\}, \quad (22b)$$

and where the *effective cost metric* \mathbf{C}_{bo} is given by (in the absence of entropic regularisation)

$$[\mathbf{C}_{bo}]_{ij} \triangleq \min_k \{ [\mathbf{C}_{ba}]_{ik} + [\mathbf{C}_{oa}]_{jk} \}. \quad (22c)$$

According to Eq. (22c), this effective cost is given by the cost of the cheapest path(s), which is intuitive. The optimal transference glued plan, \mathbf{P} , can be connected to the optimal transference plan \mathbf{P}^{bo} between \mathbf{x}^b and \mathbf{x}^o with the cost \mathbf{C}_{bo} in Eq. (22c), by marginalising on the intermediate density, i.e the W-barycentre

$$P_{ij}^{bo} = P_{ijk} \mathbf{1}_a^k. \quad (23)$$

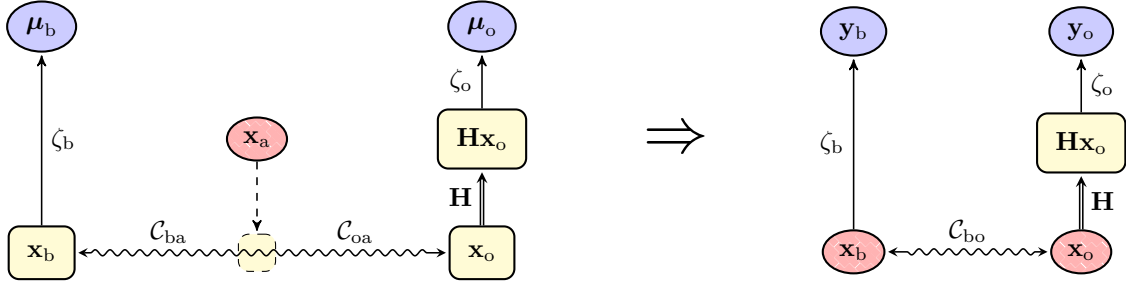


Figure 7. Trading a full hybrid OTDA problem characterised by a W-barycentre defined by the cost metrics C_{ba} and C_{oa} , with a simplified hybrid OTDA problem characterised by a single OT problem defined by an effective cost metric C_{bo} .

The solution for the analysis state x^a is given by

$$x_k^a = P_{ijk} 1_b^i 1_o^j, \quad (24)$$

by the definition of the marginals of the gluing transference plan \mathbf{P} , Eq. (19c). Yet, we do not have a direct access to the optimal
 350 gluing \mathbf{P} from the dual problem Eq. (22). This will be made simpler later on when adding the entropic regularisation to the problem.

For now, let us find an alternative solution bypassing the need for the gluing \mathbf{P} and define the map

$$\begin{aligned} \kappa^{bo} : \llbracket 1, N_b \rrbracket \times \llbracket 1, N_o \rrbracket &\mapsto \mathcal{P}(\llbracket 1, N_a \rrbracket) \\ (i, j) &\mapsto \kappa_{ij}^{bo} = \arg \min_k (C_{ba}^{ik} + C_{oa}^{jk}), \end{aligned} \quad (25)$$

355 where $\mathcal{P}(S)$ is defined as the set of all subsets of S . The set κ_{ij}^{bo} lists all the indices k that are *relays* to the transport in between the sites corresponding to index i and index j . That is why the W-barycentre can be obtained from \mathbf{P}^{bo} :

$$x_k^a = P_{ijk} 1_b^i 1_o^j = \sum_{ij} P_{ij}^{bo} \delta_{k \in \kappa_{ij}^{bo}}. \quad (26)$$

We will show in the next section how to estimate \mathbf{P}^{bo} using entropic regularisation and hence leverage Eq. (26) to compute x^a . κ^{bo} is reminiscent of the so-called *McCann interpolant* in OT theory because it is only related to the OT between x^b and x^o ,
 360 bypassing x^a , and hence the transference plan \mathbf{P}^{bo} . Please refer to Remark 7.1 by Peyré and Cuturi (2019) and to Gangbo and McCann (1996) for a description of the McCann interpolant, even when there is no Monge map. This suggests that the analysis x^a is not an interpolation of x^b and x^o in the space of values as for classical DA, but along a geodesic in a Riemannian space built on a metric derived from the Wasserstein distance.

Nonetheless, the above derivation shows that we can trade a W-barycentre problem characterised by a couple of OT problems
 365 for a single OT problem defined by an effective metric C_{bo} . This principle is schematically illustrated by Fig. 7.

This suggests a simpler two-step algorithm, where the steps consist of (i) solving a hybrid OTDA problem but with a single OT problem under an effective cost metric, which yields the analysed observables x^b and x^o and then (ii) computing the W-

barycentre of these \mathbf{x}^b and \mathbf{x}^o . Let us now go through these two steps, while adding entropic regularisation to the problem.

370 2.4.3 First step: simplified hybrid optimal transport data assimilation problem

The first step of the full OTDA algorithm is hence a simplified OTDA problem based on a single OT problem driven by the cost $C_{b|o}$. The corresponding high-level primal cost function is The associated (lower-level) primal cost function, but adding entropic regularisation ($\epsilon > 0$), is then Since \mathbf{x}^b and \mathbf{x}^o are not predetermined, the prior transference plan ν cannot be selected from $\mathcal{U}_{b|o}$ a priori. The simplest choice, which we decided to implement, is hence to set ν_{ij} to a constant, which assumes some statistical prior independence of \mathbf{x}^b and \mathbf{x}^o . A derivation of the dual problem equivalent to \mathcal{L}_ϵ can be obtained in the exact same way as in the previous subsection, although it is now less cluttered since there is only one OT to account for, instead of two. The associated Lagrangian is Again, the maps ζ_b^* and ζ_o^* are the Legendre-Fenchel transforms of the maps ζ_b and ζ_o . The variables \mathbf{f}_b and \mathbf{f}_o are Lagrange vectors; they are used to enforce the marginals of the transference plan associated to $W_{C_{b|o}}$. The unconstrained minimisation over \mathbf{P} , i.e. the inner minimisation problem in Eq., is obtained by cancelling the gradient with respect to \mathbf{P} , which yields Substituting this solution into minus the Lagrangian $-\mathcal{L}_\epsilon$ gives the regularised dual problem The notation \mathcal{J}_ϵ^* and J_ϵ^* , rather than \mathcal{L}_ϵ^* and L_ϵ^* , signifies that we work on the opposite of \mathcal{L}_ϵ^* and L_ϵ^* so as to obtain a dual problem to be minimised rather than maximised. Most importantly, we have, under conditions that will be satisfied in the following, the coincidence of the two minima $\mathcal{J}_\epsilon^* = -\mathcal{L}_\epsilon$, i.e. strong duality. Assuming one can obtain a proper correspondence between the optimal $\mathbf{f}_b, \mathbf{f}_o$ of the dual problem and $\mathbf{x}^b, \mathbf{x}^o$ of the primal problem, this implies, once again, that the primal problem can be traded for the dual problem.

Even though the regularised optimisation problem is slightly different from the unregularised one, a difference which is controlled by the value of ϵ , the new dual optimisation problem is free, i.e. without constraints. It can be solved as it is, using for instance the L-BFGS-B minimiser (Liu and Nocedal, 1989). The advantage of the regularised dual formulation is two-fold: the dual cost function is unconstrained (free optimisation) and we will trade a minimisation over $N_b \times N_o$ variables for a minimisation over $N_b + N_o$ variables. This dual formulation can be viewed as a generalised *Physical-space Statistical Analysis System* (PSAS) formalism (Courtier, 1997), an approach where classical DA algebra is mostly carried out in observation space.

Once the optimal values for \mathbf{f}_b and \mathbf{f}_o are obtained, the optimal discrete Kantorovich transference plan \mathbf{P} can be computed using Eq.. As a result, as marginals of this transference plan, the solutions for the observables are

395 2.4.3 Second step: Wasserstein barycentre

Now that we have obtained the observables \mathbf{x}^b and \mathbf{x}^o via Eq., we would like to compute their W-barycentre. The joint mass m of these observables can be computed: The high-level primal cost function of this W-barycentre problem is We have found and practised several ways to solve this problem. One way is to compute the McCann interpolant. This is theoretically elegant but Eq. did not leverage regularisation of the W-barycentre problem. Instead, the approach reported here is to use the dual optimisation problem, in conjunction with the entropic regularisation at finite $\epsilon > 0$. We leverage our knowledge of the mass

m resulting from the first step of the algorithm by enforcing the mass in the cost function, $m(\mathbf{P}) = m$. This seems redundant but it actually yields *by construction* and very naturally a numerical efficient algorithm comparable to the ad hoc log-domain scheme proposed in Sect. 4.4 of Peyré and Cuturi (2019).

Again, one way, though not the only one, to write the primal problem goes through the use of a gluing transference plan, a 3-tensor whose marginals are \mathbf{x}^b , \mathbf{x}^o , and \mathbf{x}^a . The 3-tensor ν is chosen to be $\nu_{ijk} = x_i^b x_j^o / (m N_a)$, which is uniform in k and for which $m(\nu) = m$. The resulting dual problem is This partition function is elegant but impractical since in high dimension a 3-tensor might be too large to store and compute with. However the partition function Eq. can be simplified by noticing that where we introduced the effective cost metric which is the regularised cost — known in statistics and machine learning as a *soft plus* transform — of Eq. . The 2-tensor ν_{ij} plays the same role as that of the first step of the algorithm; we choose it as $\nu_{ij} = x_i^b x_j^o / m$, for which $m(\nu) = m$. The dual problem now only involves 2-tensors and becomes numerically more efficient. Given the optimal \mathbf{f}_b and \mathbf{f}_o , the (glued) optimal transference plan \mathbf{P}^{boa} is formally given by The W -barycentre \mathbf{x}^a is then given as a marginal of \mathbf{P}^{boa} . Because of the normalisation of the transference plan to m , the entropic regularisation exhibits a $\varepsilon m \ln Z_\varepsilon$ instead of $\varepsilon Z_\varepsilon$. This systematically enforces normalisation in the computations of the gradients, as well as in the course of the numerical optimisation of the dual cost function, de facto working in log-domain. We experienced more stable computations and the ability to reach smaller ε , as compared to the case without normalisation. This completes the solution through the 2-step OTDA algorithm To avoid making a too large detour, the derivation of this algorithm is presented in Appendix B.

2.4.3 Classical data assimilation as a particular case

The primal problem Eq. (17) of classical DA reads

$$\mathcal{L}_{\text{cl}} = \min_{\mathbf{x}^a \in \mathcal{O}_{N_a}^+} \{ \zeta_b(\mathbf{y}^b - \mathbf{x}^a) + \zeta_o(\mathbf{y}^o - \mathbf{H}\mathbf{x}^a) \}. \quad (27)$$

Let us see how the OTDA formalism Eq. (22) can account for classical DA. In the context of classical DA, the observable spaces for \mathbf{x}^b , \mathbf{x}^o , and \mathbf{x}^a are assumed to coincide by construction. Let us then define the cost matrices

$$[\mathbf{C}_{\text{ba}}^\infty]_{ij} \triangleq [\mathbf{C}_{\text{oa}}^\infty]_{ij} \triangleq \begin{cases} 0 & \text{if } i = j \\ +\infty & \text{if } i \neq j \end{cases}, \quad (28)$$

i.e. it is assumed that the cost of moving masses is as large as can be. Looking back at Eq. (19) but with these costs, it is clear that in order to avoid the primal cost function to be $+\infty$, the transference plan P_{ijk} must always be 0 unless $i = j = k$. But this implies from the definition of \mathcal{U}_{boa} that the observables coincide: $\mathbf{x}^b = \mathbf{x}^o = \mathbf{x}^a$ and that their mass is given by $m(\mathbf{P})$. Hence in this limit where the specific cost matrices are equal to $\mathbf{C}_{\text{ba}}^\infty$ and $\mathbf{C}_{\text{oa}}^\infty$, the OTDA primal problem is mathematically equivalent to the classical DA primal problem. And classical DA can be seen as become mathematically equivalent to classical DA. Hence, classical DA is a limit case of the hybrid OTDA problem with the specific cost matrices $\mathbf{C}_{\text{ba}}^\infty$ and $\mathbf{C}_{\text{oa}}^\infty$ OTDA. Note that from its definition Eq. (22c), the effective cost \mathbf{C}_{bo} obtained from $\mathbf{C}_{\text{ba}}^\infty$ and $\mathbf{C}_{\text{oa}}^\infty$ coincide with $\mathbf{C}_{\text{bo}}^\infty \triangleq \mathbf{C}_{\text{ba}}^\infty = \mathbf{C}_{\text{oa}}^\infty$.

2.5 A direct algorithmic solution

The two-step approach of Sect. 2.4 has the merit to connect to the traditional W-barycentre problem, by first estimating \mathbf{x}^b and \mathbf{x}^o , and later computing the W-barycentre in between both states. It also suggests the existence of the effective cost metric of the problem. However, going through its consecutive steps may not be necessary for pure computational purposes. Here we
 435 describe a direct approach that yields the analysis of the OTDA problem. It is less enlightening but is practical and will be used in the subsequent illustrations of the present paper.

An alternative formulation to the primal problem Eq. (19) relies on two transference plans \mathbf{P}^{ba} and \mathbf{P}^{oa} corresponding to the two transports of the underlying W-barycentre problem, instead of the gluing one. Moreover, entropic regularisation is enforced via $\mathcal{K}(\mathbf{P}^{ba}|\boldsymbol{\nu}^{ba})$ and $\mathcal{K}(\mathbf{P}^{oa}|\boldsymbol{\nu}^{oa})$. The corresponding optimisation problem reads

$$440 \quad \mathcal{L} = \min_{\substack{\mathbf{x}^b \in \mathcal{O}_b^+ \\ \mathbf{x}^o \in \mathcal{O}_o^+ \\ \mathbf{x}^a \in \mathcal{O}_a^+}} [\zeta_b(\mathbf{y}^b - \mathbf{x}^b) + \zeta_o(\mathbf{y}^o - \mathbf{H}\mathbf{x}^o) \\ + \min_{\mathbf{P}^{ba} \in \mathcal{U}_{ba} \mathbf{P}^{oa} \in \mathcal{U}_{oa}} \{ \varepsilon \mathcal{K}(\mathbf{P}^{ba}|\boldsymbol{\nu}^{ba}) + \varepsilon \mathcal{K}(\mathbf{P}^{oa}|\boldsymbol{\nu}^{oa}) + P_{ik}^{ba} C_{ba}^{ik} + P_{jk}^{oa} C_{oa}^{jk} \}], \quad (29a)$$

where the admissible sets of transference plans \mathbf{P}^{ba} and \mathbf{P}^{oa} are defined by

$$\mathcal{U}_{ba} \triangleq \{ \mathbf{P} \in \mathcal{O}_{b,a}^+ : \mathbf{P}\mathbf{1}_a = \mathbf{x}^b, \quad \mathbf{P}^\top \mathbf{1}_b = \mathbf{x}^a \}, \quad (29b)$$

$$\mathcal{U}_{oa} \triangleq \{ \mathbf{P} \in \mathcal{O}_{o,a}^+ : \mathbf{P}\mathbf{1}_a = \mathbf{x}^o, \quad \mathbf{P}^\top \mathbf{1}_b = \mathbf{x}^a \}. \quad (29c)$$

445 Following the same type of derivation as reported in the previous sections [and Appendix B](#), the corresponding dual problem to be minimised is obtained as

$$J_\varepsilon^* = \min_{\mathbf{f}_b \in \mathbb{R}^{N_b} \mathbf{f}_o \in \mathbb{R}^{N_o} \mathbf{f}_a \in \mathbb{R}^{N_a}} J_\varepsilon^*(\mathbf{f}_b, \mathbf{f}_o, \mathbf{f}_a), \quad (30a)$$

where, discarding the constant $-\varepsilon m(\boldsymbol{\nu}^{ba}) - \varepsilon m(\boldsymbol{\nu}^{oa})$, the associated regularised Lagrangian is

$$J_\varepsilon^*(\mathbf{f}_b, \mathbf{f}_o, \mathbf{f}_a) = \varepsilon Z_\varepsilon^{ba}(\mathbf{f}_b, \mathbf{f}_a) + \varepsilon Z_\varepsilon^{oa}(\mathbf{f}_o, \mathbf{f}_a) + \zeta_b^*(\mathbf{f}_b) + \zeta_o^*(\mathbf{f}_o) - \mathbf{f}_b^\top \mathbf{y}^b - \mathbf{f}_o^\top \mathbf{y}^o, \quad (30b)$$

450 with a *partition function* associated to each transport:

$$Z_\varepsilon^{ba} \triangleq \sum_{ik} P_{ik}^{ba}, \quad Z_\varepsilon^{oa} \triangleq \sum_{jk} P_{jk}^{oa} \quad (30c)$$

where

$$P_{ik}^{ba} = \nu_{ik}^{ba} e^{(f_b^i + f_a^k - C_{ba}^{ik})/\varepsilon}, \quad P_{jk}^{oa} = \nu_{jk}^{oa} e^{(f_o^j + f_a^k - C_{oa}^{jk})/\varepsilon}. \quad (30d)$$

It turns out that the optimal \mathbf{f}_a can be obtained analytically as a function of \mathbf{f}_b and \mathbf{f}_o , which we checked makes the optimisation
 455 numerically more efficient and robust. Indeed, let us introduce $\psi_k \triangleq e^{f_a^k/\varepsilon}$. We could optimise $J_\varepsilon^*(\mathbf{f}_b, \mathbf{f}_o, \mathbf{f}_a = \varepsilon \ln \boldsymbol{\psi})$ on $\boldsymbol{\psi}$:

$$0 = \partial_{\psi_k} J_\varepsilon^*(\mathbf{f}_b, \mathbf{f}_o, \mathbf{f}_a) = \sum_i \nu_{ik}^{ba} e^{(f_b^i - C_{ba}^{ik})/\varepsilon} - \frac{1}{\psi_k^2} \sum_j \nu_{jk}^{oa} e^{(f_o^j - C_{oa}^{jk})/\varepsilon}, \quad (31)$$

yielding the solution

$$\psi_k^2 = \frac{Z_{\varepsilon,k}^{\text{oa}}}{Z_{\varepsilon,k}^{\text{ba}}}, \quad Z_{\varepsilon,k}^{\text{oa}} \triangleq \sum_j \nu_{jk}^{\text{oa}} e^{(f_o^l H_l^j - C_{\text{oa}}^{jk})/\varepsilon}, \quad Z_{\varepsilon,k}^{\text{ba}} \triangleq \sum_i \nu_{ik}^{\text{ba}} e^{(f_b^i - C_{\text{ba}}^{ik})/\varepsilon}. \quad (32)$$

Up to irrelevant constants, the resulting effective cost function using the optimal ψ_k is

$$460 \quad J_\varepsilon^*(\mathbf{f}_b, \mathbf{f}_o) = 2\varepsilon \sum_k \sqrt{Z_{\varepsilon,k}^{\text{ba}} Z_{\varepsilon,k}^{\text{oa}}} + \zeta_b^*(\mathbf{f}_b) + \zeta_o^*(\mathbf{f}_o) - \mathbf{f}_b^\top \mathbf{y}^b - \mathbf{f}_o^\top \mathbf{y}^o. \quad (33)$$

Now, the optimal W-barycentre \mathbf{x}^a is given by either $x_k^a = P_{ik}^{\text{ba}} 1_b^i$ or $x_k^a = P_{jk}^{\text{oa}} 1_o^j$, i.e.

$$x_k^a = \psi_k Z_{\varepsilon,k}^{\text{ba}} = \frac{1}{\psi_k} Z_{\varepsilon,k}^{\text{oa}}, \quad (34)$$

from which we can infer the ψ_k -free expression

$$x_k^a = \sqrt{Z_{\varepsilon,k}^{\text{ba}} Z_{\varepsilon,k}^{\text{oa}}}. \quad (35)$$

465 It is also useful to retrieve the optimal value of \mathbf{f}_a and obtain

$$f_a^k = \varepsilon \ln \psi_k = \frac{\varepsilon}{2} \ln \left(\frac{Z_{\varepsilon,k}^{\text{oa}}}{Z_{\varepsilon,k}^{\text{ba}}} \right), \quad (36)$$

so that we can compute the other two analysed observables, \mathbf{x}^b and \mathbf{x}^o , using

$$x_i^b = P_{ik}^{\text{ba}} 1_a^k = \sum_k \psi_k \nu_{ik}^{\text{ba}} e^{(f_b^i - C_{\text{ba}}^{ik})/\varepsilon} = e^{f_b^i/\varepsilon} \sum_k \nu_{ik}^{\text{ba}} e^{(f_a^k - C_{\text{ba}}^{ik})/\varepsilon}, \quad (37a)$$

$$x_j^o = P_{jk}^{\text{oa}} 1_a^k = \sum_k \frac{1}{\psi_k} \nu_{jk}^{\text{oa}} e^{(f_o^l H_l^j - C_{\text{oa}}^{jk})/\varepsilon} = e^{(f_o^l H_l^j)/\varepsilon} \sum_k \nu_{jk}^{\text{oa}} e^{(-f_a^k - C_{\text{oa}}^{jk})/\varepsilon}. \quad (37b)$$

470 Note that most of these expressions can be assessed in a robust way in the log-domain. For instance we use in practice, equivalently to Eqs. (35,37):

$$\varepsilon \ln x_k^a = \frac{\varepsilon}{2} \ln \sum_i \nu_{ik}^{\text{ba}} e^{(f_b^i - C_{\text{ba}}^{ik})/\varepsilon} + \frac{\varepsilon}{2} \ln \sum_j \nu_{jk}^{\text{oa}} e^{(f_o^l H_l^j - C_{\text{oa}}^{jk})/\varepsilon}, \quad (38a)$$

$$\varepsilon \ln x_i^b = f_b^i + \varepsilon \ln \sum_k \nu_{ik}^{\text{ba}} e^{(f_a^k - C_{\text{ba}}^{ik})/\varepsilon}, \quad (38b)$$

$$\varepsilon \ln x_j^o = f_o^l H_l^j + \varepsilon \ln \sum_k \nu_{jk}^{\text{oa}} e^{(-f_a^k - C_{\text{oa}}^{jk})/\varepsilon}. \quad (38c)$$

475 3 Numerical illustrations

In this section, we showcase a selection of OTDA ~~3D-Var~~ 3D-Var analyses. These are meant to stress the versatility of the formalism and the diverse solutions it offers, with significantly more degrees of freedom than in classical DA. The OTDA state analysis is carried out using Sect. 2.5 and its formulas. Unless specifically discussed, entropic regularisation is used with $\varepsilon = 10^{-3}$. The dual cost function Eq. (33) is minimised using the quasi-Newton method L-BFGS-B (Liu and Nocedal, 1989),

480 which yields the optimal $\mathbf{f}_b, \mathbf{f}_o$. Then Eq. (38) is employed to compute $\mathbf{x}^b, \mathbf{x}^o$, and \mathbf{x}^a .

3.1 One-dimensional examples

Considering the case where the physical space of the fields is one-dimensional, we build bell-shaped observations \mathbf{y}^b and \mathbf{y}^o , related to an observable space of size $N_b = N_o = N_a = 10^2$ shared by \mathbf{x}^b , \mathbf{x}^o and \mathbf{x}^a . Since \mathbf{y}^b is a fully observed instance of \mathbf{x}^b , we have $\mathfrak{N}_b = N_b = 10^2$, while \mathfrak{N}_o may differ from N_o depending on the definition of the observation operator \mathbf{H} . We
 485 choose (Gaussian statistics)

$$\zeta_b(\mathbf{e}_b) = \frac{1}{2\sigma_b^2} \|\mathbf{e}_b\|^2, \quad \zeta_b(\mathbf{e}_o) = \frac{1}{2\sigma_o^2} \|\mathbf{e}_o\|^2, \quad (39)$$

with $\sigma_b = \sigma_o = 10^{-2}$. The states are discretised over the interval $[0, 1]$ at sites/grid cells $r_i^\star = (i - \frac{1}{2})/N_\star$ for $i \in \llbracket 1, N_\star \rrbracket$, with $\star = b, o, a$. Unless otherwise specified, the cost metric has a quadratic dependence with the distance between sites, i.e. $[\mathbf{C}_{ba}]_{ik} = |r_i^b - r_k^a|^2$ and $[\mathbf{C}_{oa}]_{jk} = |r_j^o - r_k^a|^2$. This is our reference setup. The observation operator and the mass of the
 490 observations \mathbf{y}^b and \mathbf{y}^o will be described specified for each experiment.

We consider four experiments where we choose to vary key parameters in the OTDA setup.

3.1.1 Varying the imbalance of the observation states

In the first experiment, the system is fully observed with $\mathbf{H} = \mathbf{I}$. We choose $m(\mathbf{y}^b) = 1$ and the mass of \mathbf{y}^o to be in the set $m(\mathbf{y}^o) \in \{0.5, 1, 1.5\}$, all the other parameters being fixed to the reference. The results are displayed in Fig. 8. Panel (a)
 495 corresponds to the case $m(\mathbf{y}^o) = 0.5$. The resulting mass of the analysed observables is then $m(\mathbf{x}^a) = m(\mathbf{x}^b) = m(\mathbf{x}^o) = 0.79$. The adjustment of \mathbf{x}^b compared to \mathbf{y}^b , and the adjustment of \mathbf{x}^o compared to \mathbf{y}^o , which are required to balance \mathbf{x}^b , \mathbf{x}^o are patent. Panel (b) corresponds to the case $m(\mathbf{y}^o) = 1$. The resulting mass of the analysed observables is then $m(\mathbf{x}^a) = m(\mathbf{x}^b) = m(\mathbf{x}^o) = 1$. No adjustment is required here since $m(\mathbf{y}^o) = m(\mathbf{y}^b)$, and \mathbf{x}^o and \mathbf{y}^o , as well as \mathbf{x}^b and \mathbf{y}^b coincide. Finally, the mass of \mathbf{y}^o is set to $m(\mathbf{y}^o) = 1.5$ in panel (c). The resulting mass of the analysed observables is then $m(\mathbf{x}^a) =$
 500 $m(\mathbf{x}^b) = m(\mathbf{x}^o) = 1.34$. The adjustment of \mathbf{x}^b compared to \mathbf{y}^b , and the adjustment of \mathbf{x}^o compared to \mathbf{y}^o , which are required to balance \mathbf{x}^b , \mathbf{x}^o are visually obvious, but the balancing goes in the opposite direction compared to panel (a), as expected.

3.1.2 Varying the sparseness of the observation operator

In this second experiment, all the other parameters being fixed to their reference value, only a fraction of the domain is observed, over $[0, \frac{1}{4}]$, $[0, \frac{1}{2}]$ and $[0, \frac{3}{4}]$, where $\mathbf{H} \in \mathcal{O}_{\mathfrak{N}_o}^+ \times N_o$ with $\mathfrak{N}_o = N_o/4, N_o/2, 3N_o/4$, and $H_l^j = \delta_{l,j}$ for $l \in \llbracket 1, \mathfrak{N}_o \rrbracket$
 505 and $j \in \llbracket 1, N_o \rrbracket$.

The masses of the states that are built to generate \mathbf{y}^b and \mathbf{y}^o , before applying any observation operator, are set to 1 and 1.5, respectively. As a result, we have $m(\mathbf{y}^b) = 1$ but $m(\mathbf{y}^o)$ may depart from 1.5 depending on \mathbf{H} . The fully observed case corresponds to panel (c) of Fig. 8. The results are displayed in Fig. 9. It shows how smooth the OTDA solution can be compared to that of classical DA. Yet, as in panel (a), OTDA can also handle obviously diverging sources of information as in the case
 510 where the support of \mathbf{H} is $[0, \frac{1}{4}]$ and where \mathbf{y}^o and \mathbf{y}^b can be seen are barely consistent. In that case, the OTDA solution is smooth but bimodal.

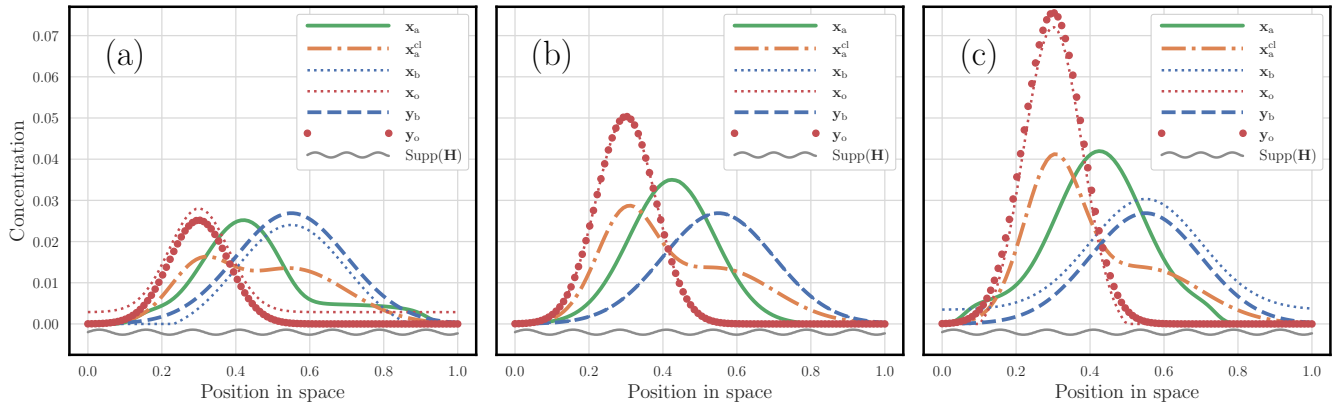


Figure 8.

A hybrid OTDA 3D-Var-3D-Var analysis with one-dimensional physical states, where only the mass of \mathbf{y}° is varied. Its mass is $m(\mathbf{y}^\circ) = 0.5$ in panel (a), $m(\mathbf{y}^\circ) = 1$ in panel (b), and $m(\mathbf{y}^\circ) = 1.5$ in panel (c). The dashed blue curve corresponds to the first guess \mathbf{y}^b , the red dots correspond to the observations \mathbf{y}° , the analysis state \mathbf{x}^a is the solid green curve, the analysed observables \mathbf{x}^b and \mathbf{x}° are blue and red dotted curves, respectively. The support of the observation is underlined by a wavy grey segment. The corresponding classical analysis is also plotted with a dashed-dotted orange curve. The x-axis corresponds to the position in space; the y-axis corresponds to the concentration value of the fields.

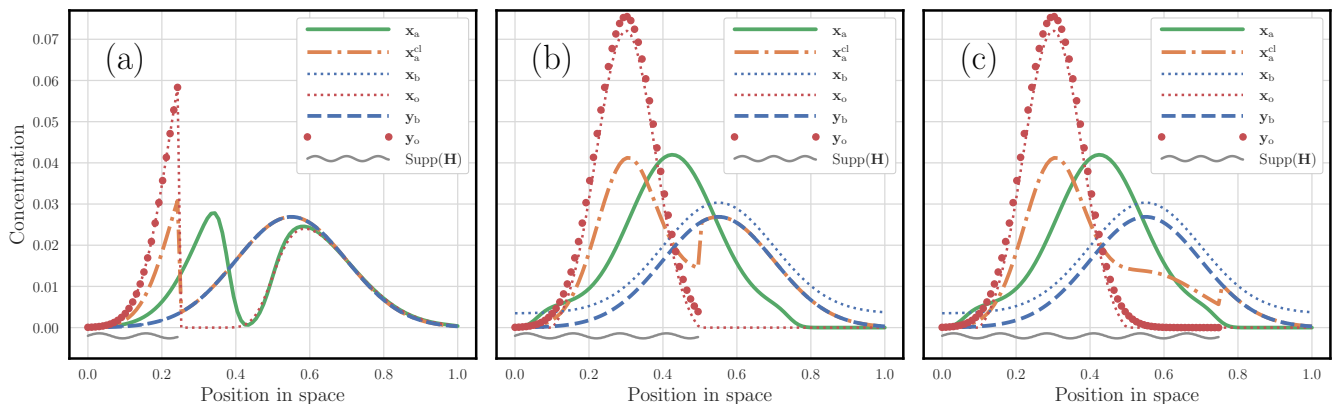


Figure 9.

A hybrid OTDA 3D-Var-3D-Var analysis with one-dimensional physical states, where the observation operator is increasingly sparse. The support of \mathbf{H} is $[0, \frac{1}{4}]$ for panel (a), $[0, \frac{1}{2}]$ for panel (b), and $[0, \frac{3}{4}]$ for panel (c). See Fig. 8 for the description of the legend.

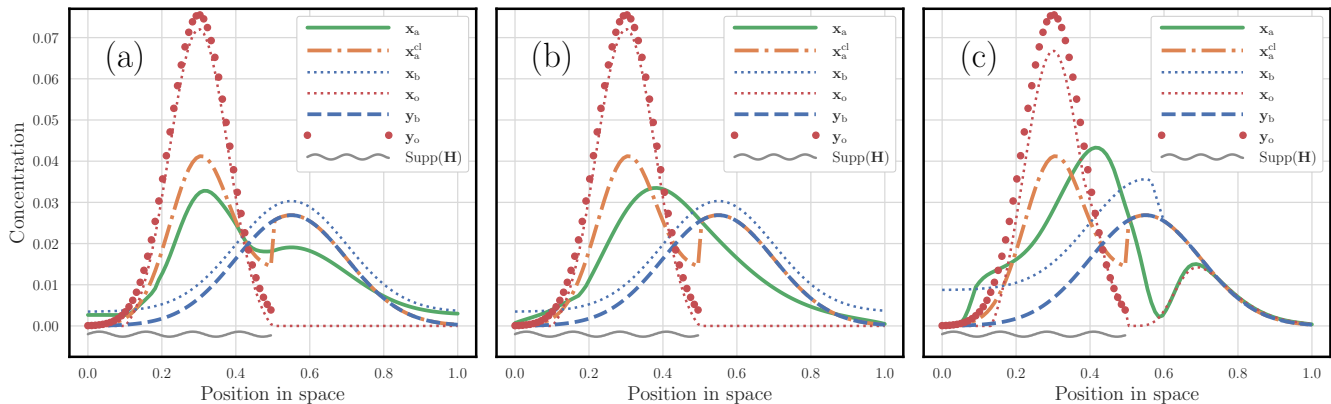


Figure 10.

A hybrid OTDA ~~3D-Var~~ ~~3D-Var~~ analysis with one-dimensional physical states, where the cost metrics are changed. See the bulk of the text for a definition of those three cost metrics. See Fig. 8 for the description of the legend.

3.1.3 Changing the nature of the cost metric

In this third experiment, we choose the cost metric to be of the form $[\mathbf{C}_{ba}]_{ik} = |r_i^b - r_k^a|^\alpha$ and $[\mathbf{C}_{oa}]_{jk} = |r_j^o - r_k^a|^\alpha$. Only half of the domain is observed over $[0, \frac{1}{2}]$, as in the case of Fig. 9, panel (b). Since the mass of the state used to produce y^o is 1.5, we have a slightly different $m(y^o) = 1.49$, the rest of the mass being located in the unobserved part of the domain. All of the other parameters follow the reference setup. The results are displayed in Fig. 10. For panel (a), α is set to 0.5. For panel (b), α is set to 1. For panel (c), the cost metric is piecewise; it is quadratic, i.e. $\alpha = 2$, for pairs of sites separated by less than 10^{-1} , i.e. $|r_i^b - r_k^a| = |r_j^o - r_k^a| \leq 10^{-1}$, whereas for pairs of sites beyond this range, the costs are chosen to be infinite. Hence transport is prohibited beyond a distance of 10^{-1} . The case of a pure quadratic cost corresponds to panel (b) of Fig. 9. The impact on the shape of the OTDA analysis is very significant, and suggests that one could easily tailor their own cost to suit their specific DA problem.

3.1.4 Classical data assimilation as a subcase of the hybrid optimal transport data assimilation

In the fourth experiment, we would like to numerically check the theoretical prediction of Sect. 2.4.3. Consider again the reference configuration. But only half of the domain, over $[0, \frac{1}{2}]$, is observed, $\mathbf{H} \in \mathcal{O}_{\mathfrak{N}_o \times N_o}^+$ with $\mathfrak{N}_o = N_o/2$ and $H_l^j = \delta_{l,j}$ for $l \in \llbracket 1, \mathfrak{N}_o \rrbracket$ and $j \in \llbracket 1, N_o \rrbracket$. Most importantly, the cost metric has a quadratic dependence with the distance between sites, i.e. $[\mathbf{C}_{ba}]_{ik} = \lambda |r_i^b - r_k^a|^2$ and $[\mathbf{C}_{oa}]_{jk} = \lambda |r_j^o - r_k^a|^2$. The case $\lambda = 1$ corresponds to panel (b) of Fig. 9. Figure 11 shows the results corresponding to: $\lambda = 10^3$ for panel (a), $\lambda = 10^4$ for panel (b), and $\lambda = 10^6$ for panel (c). When λ is increased, the OTDA analysis should tend to the classical DA solution. This is indeed corroborated by Fig. 11 and comforts the claim of

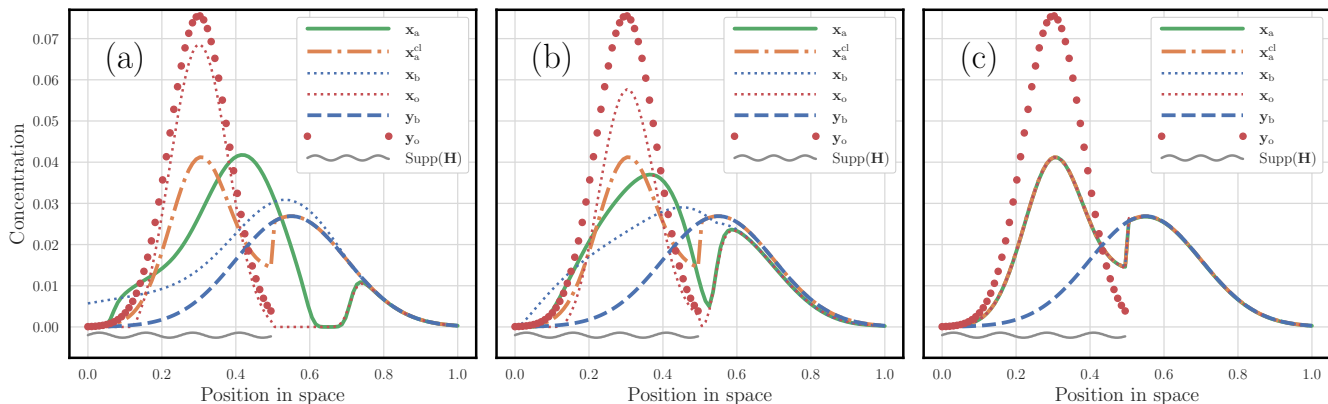


Figure 11. Scaling up the cost metrics λC_{ba} and λC_{oa} with increasing λ , the OTDA analysis converges to the classical DA analysis. Panels (a), (b), and (c) correspond to the scaling values $\lambda = 10^3, 10^4, 10^6$, respectively. See Fig. 8 for the description of the legend.

Sect. 2.4.3. Note that, as opposed to the three earlier experiments, we had here to tune ε since the wide range of λ has a
 530 significant impact in the balance of the key terms of the cost function (transport cost, discrepancy errors, and regularisation).

3.2 Two-dimensional examples

Considering the case where the physical space of the fields is two-dimensional, we perform a couple of **3D-Var-3D-Var**
 analysis on concentration fields (puffs of a pollutant). The states are discretised in the domain $[0, 1]^2$ at sites/grid cells $\mathbf{r}_{i,j}^* =$
 $((i - \frac{1}{2})/N_x^*, (j - \frac{1}{2})/N_y^*)$ for $(i, j) \in \llbracket 1, N_x^* \rrbracket \times \llbracket 1, N_y^* \rrbracket$, with $\star = b, o, a$. We choose $N_b^x = N_b^y = N_o^x = N_o^y = N_a^x = N_a^y =$
 535 10^2 , such that $N_b = N_o = N_a = 10^4$. Hence, the number of control variables is 3×10^4 . The observation vectors are \mathbf{y}^b and
 \mathbf{y}^o . Since \mathbf{y}^b is a fully observed instance of \mathbf{x}^b , we have $\mathfrak{N}_b = N_b$, while \mathfrak{N}_o may differ from N_o depending on the definition
 of the observation operator \mathbf{H} . Moreover, we choose (Gaussian statistics)

$$\zeta_b(\mathbf{e}_b) = \frac{1}{2\sigma_b^2} \|\mathbf{e}_b\|^2, \quad \zeta_b(\mathbf{e}_o) = \frac{1}{2\sigma_o^2} \|\mathbf{e}_o\|^2, \quad (40)$$

with $\sigma_b = \sigma_o = 10^{-2}$. The entropic regularisation parameter is set to $\varepsilon = 10^{-3}$.

540 The first analysis is displayed in Fig. 12. The observation operator \mathbf{H} is the identity but its support is restricted to the
 subdomain $[0, 0.6]^2$. The plumes of pollutants \mathbf{y}^b and \mathbf{y}^o are generated from states formed as combinations of bell-like puffs.
 The system is unbalanced with $m(\mathbf{y}^b) = 1.35$, $m(\mathbf{y}^o) = 0.73$. The cost metric has a quadratic dependence with the distance
 between sites, i.e. $[C_{ba}]_{ik} = \|\mathbf{r}_i^b - \mathbf{r}_k^a\|_2^2$ and $[C_{oa}]_{jk} = \|\mathbf{r}_j^o - \mathbf{r}_k^a\|_2^2$. The OTDA analysis is clearly smoother than the classical
 solution. The classical solution does not cope very well with the seemingly disagreeing sources of information \mathbf{y}^b and \mathbf{y}^o ,
 545 which generates sharp transitions in the classical analysis. If \mathbf{y}^b and \mathbf{y}^o were consistently obtained from a truth perturbed with
 errors with short-range correlation, then the classical analysis would be as good as can be, while the OTDA solution may be

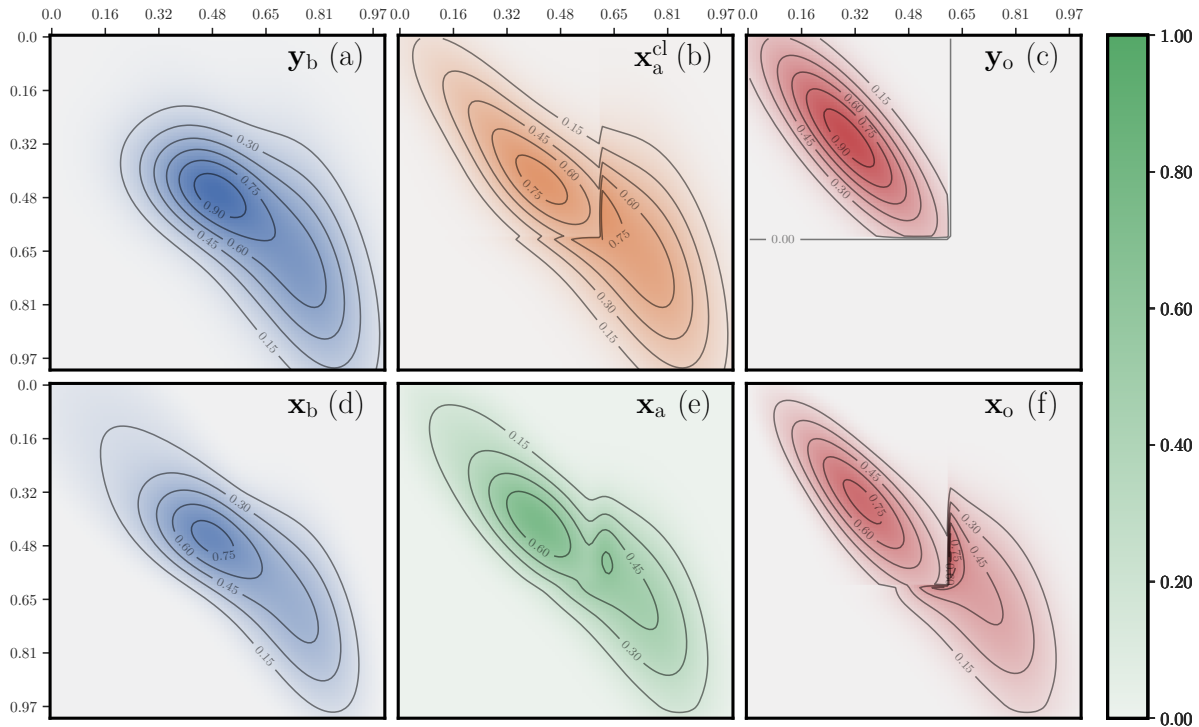


Figure 12. Two-dimensions concentration maps (plumes) of a hybrid OTDA analysis for the first configuration. The observations \mathbf{y}^b and \mathbf{y}^o , the analysed observables \mathbf{x}^b , \mathbf{x}^a i.e. the state analysis, \mathbf{x}^o , and the corresponding classical DA analysis \mathbf{x}_{cl}^a are displayed. All fields are rescaled so their joint maximum is 1. All heatmaps use the same scale. The colour bar represents a unified contrast scale for the diverse field concentrations.

too safe. However, if one believes that structural errors and in particular location errors can impact \mathbf{y}^b and \mathbf{y}^o , then the classical solution is improper and the OTDA analysis preferable.

The second analysis is displayed in Fig. 13. The support of the observation operator \mathbf{H} is again contained within the subdomain $[0, 0.6]^2$ but only one of four grid cells are actually observed in this area. The observation states \mathbf{y}^b and \mathbf{y}^o are generated from the same states as for Fig. 12. The system is unbalanced with $m(\mathbf{y}^b) = 1.35$, $m(\mathbf{y}^o) = 0.18$. The cost metric is defined to be the same as in Fig. 12. The OTDA analysis is even smoother in this case as compared to the classical DA analysis. It is much less impacted by the sparseness of the observation operator. The classical solution has to account for the staggered observations in the top left corner of the domain because the first guess in that region is weak. By contrast, the OTDA solution
555 assumes that location errors are possible and it hence moves around the mass corresponding to these observations, so that the structure of the observation operator is not as impactful on the analysis.

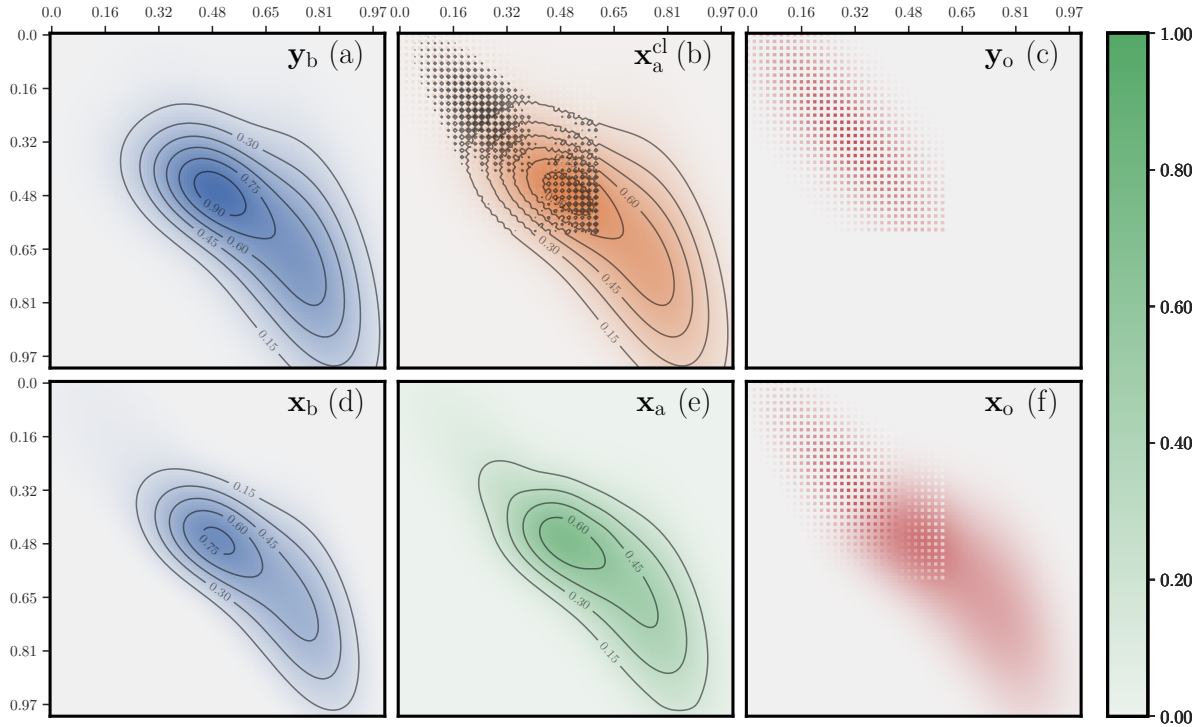


Figure 13. Two-dimensions concentrations maps of a hybrid OTDA analysis for the second configuration. The observations \mathbf{y}^b and \mathbf{y}^o , the analysed observables \mathbf{x}^b , \mathbf{x}^a i.e. the state analysis, \mathbf{x}^o , and the corresponding classical DA analysis. Compared to Fig. 12, only \mathbf{H} has changed. The level sets in panels (c) and (f) are omitted since they are driven by the staggered observation operator.

4 Uncertainty quantification

In this section, we compute the posterior error covariance matrix \mathbf{P}^a associated to the state analysis \mathbf{x}^a , in order to complete the OTDA 3D-Var-3D-Var analysis description. There are many ways to proceed depending on the chosen regularisation and on the targeted degree of generality. Here, for the sake of consistency, we report on the way to derive \mathbf{P}^a following the computation of the analysis state \mathbf{x}^a proposed in Sect. 2.5.

4.1 Mathematical results

Let us denote the compounded vectors of the observations, of the Lagrange multipliers, and of the observables, as well as the compounded observation operator by

$$565 \quad \mathbf{y} \triangleq \begin{bmatrix} \mathbf{y}^b \\ \mathbf{y}^o \end{bmatrix}, \quad \mathbf{f} \triangleq \begin{bmatrix} \mathbf{f}_b \\ \mathbf{f}_o \end{bmatrix}, \quad \mathbf{x} \triangleq \begin{bmatrix} \mathbf{x}^b \\ \mathbf{x}^o \end{bmatrix}, \quad \mathcal{H} \triangleq \begin{bmatrix} \mathbf{I}_b & \mathbf{0} \\ \mathbf{0} & \mathbf{H} \end{bmatrix}, \quad (41)$$

of size $\mathfrak{N}_b + \mathfrak{N}_o$, $\mathfrak{N}_b + \mathfrak{N}_o$, $N_b + N_o$, and $(\mathfrak{N}_b + \mathfrak{N}_o) \times (N_b + N_o)$, respectively. Similarly, we define the sum of the error statistics by $\zeta(\mathbf{f}) \triangleq \zeta_b(\mathbf{f}_b) + \zeta_o(\mathbf{f}_o)$, whose Legendre-Fenchel transform is $\zeta^*(\mathbf{f}) = \zeta_b^*(\mathbf{f}_b) + \zeta_o^*(\mathbf{f}_o)$. Using this notation, we can recapitulate the key results of Sect. 2.5: the effective dual cost function is

$$J_\varepsilon^*(\mathbf{f}) \triangleq \varepsilon Z_\varepsilon(\mathbf{f}) + \zeta^*(\mathbf{f}) - \mathbf{f}^\top \mathbf{y}, \quad Z_\varepsilon(\mathbf{f}) \triangleq 2 \sum_k \sqrt{Z_{\varepsilon,k}^{\text{ba}}(\mathbf{f}_b) Z_{\varepsilon,k}^{\text{oa}}(\mathbf{f}_o)}, \quad (42)$$

570 and the analysis state reads

$$x_k^a(\mathbf{f}) = \sqrt{Z_{\varepsilon,k}^{\text{ba}}(\mathbf{f}_b) Z_{\varepsilon,k}^{\text{oa}}(\mathbf{f}_o)}, \quad (43)$$

where the dependence of the analysis state and the partition functions on \mathbf{f} , \mathbf{f}_b and \mathbf{f}_o is now emphasised and made explicit.

Any prior source of error in the system stems from [the information vector \$\mathbf{y}\$](#) , and hence drives the posterior error in the analysis \mathbf{x}^a . That is why we are interested in the sensitivity of \mathbf{x}^a with respect to \mathbf{y} , i.e. $\delta \mathbf{x}^a = \partial_{\mathbf{y}} \mathbf{x}^a \delta \mathbf{y}$. Denoting the expectation

575 operator by \mathbb{E} , the error covariance matrix is then defined by

$$\mathbf{P}^a = \mathbb{E} \left[\delta \mathbf{x}^a (\delta \mathbf{x}^a)^\top \right] = (\partial_{\mathbf{y}} \mathbf{x}^a) \mathbb{E} [\delta \mathbf{y} \delta \mathbf{y}^\top] (\partial_{\mathbf{y}} \mathbf{x}^a)^\top = (\partial_{\mathbf{y}} \mathbf{x}^a) (\partial_{\mathbf{f}}^2 \zeta^*) (\partial_{\mathbf{y}} \mathbf{x}^a)^\top, \quad (44)$$

from which a matrix factor \mathbf{X}^a of \mathbf{P}^a , i.e. which satisfies $\mathbf{P}^a = \mathbf{X}^a (\mathbf{X}^a)^\top$ and whose expressions are usually much shorter than those of \mathbf{P}^a , can be extracted, up to the multiplication by an orthogonal matrix on the right:

$$\mathbf{X}^a = \partial_{\mathbf{y}} \mathbf{x}^a (\partial_{\mathbf{f}}^2 \zeta^*)^{\frac{1}{2}}. \quad (45)$$

580 To compute the sensitivity matrix $\partial_{\mathbf{y}} \mathbf{x}^a$, we leverage the stationarity of the dual cost function at the minimum:

$$\partial_{\mathbf{f}} J_\varepsilon^*(\mathbf{f}(\mathbf{y}), \mathbf{y}) = \mathbf{0}, \quad (46)$$

and resort to the implicit function theorem:

$$\mathbf{0} = d_{\mathbf{y}} \partial_{\mathbf{f}} J_\varepsilon^*(\mathbf{f}(\mathbf{y}), \mathbf{y}) = \partial_{\mathbf{f}}^2 J_\varepsilon^* \partial_{\mathbf{y}} \mathbf{f} + \partial_{\mathbf{f}} \partial_{\mathbf{y}} J_\varepsilon^*, \quad (47)$$

which yields

$$585 \quad \partial_{\mathbf{y}} \mathbf{f} = - [\partial_{\mathbf{f}}^2 J_\varepsilon^*]^{-1} \partial_{\mathbf{f}} \partial_{\mathbf{y}} J_\varepsilon^* = [\partial_{\mathbf{f}}^2 J_\varepsilon^*]^{-1}, \quad (48)$$

since $\partial_{\mathbf{f}} \partial_{\mathbf{y}} J_\varepsilon^* = -\mathbf{I}_{b_o}$, where \mathbf{I}_{b_o} is the identity matrix in the compounded observation space $\mathbb{R}^{\mathfrak{N}_b + \mathfrak{N}_o}$. The sensitivity $\partial_{\mathbf{y}} \mathbf{x}^a$ can now be computed using the Leibniz chain rule and Eq. (48):

$$\frac{\partial \mathbf{x}^a}{\partial \mathbf{y}} = \frac{\partial \mathbf{x}^a}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \mathbf{y}} = \partial_{\mathbf{f}} \mathbf{x}^a [\partial_{\mathbf{f}}^2 J_\varepsilon^*]^{-1}. \quad (49)$$

Let us now compute the Jacobian and Hessian in the right-hand side of Eq. (49). To that end and in order to externalise the

590 observation operator, we introduce \hat{Z}_ε and $\hat{\mathbf{x}}_a$ such that

$$\hat{Z}_\varepsilon(\boldsymbol{\eta} = \mathcal{H}^\top \mathbf{f}) \triangleq Z_\varepsilon(\mathbf{f}), \quad \hat{\mathbf{x}}_a(\boldsymbol{\eta} = \mathcal{H}^\top \mathbf{f}) \triangleq \mathbf{x}^a(\mathbf{f}), \quad (50)$$

and the related Jacobian and Hessians

$$\mathbf{\Omega}_{\text{bo},\text{a}} \triangleq \partial_{\boldsymbol{\eta}} \hat{\mathbf{x}}_{\text{a}}, \quad \mathbf{\Omega}_{\text{bo},\text{bo}} \triangleq \varepsilon \partial_{\boldsymbol{\eta}}^2 \hat{Z}_{\varepsilon}, \quad \mathbf{\Lambda}_{\text{bo},\text{bo}} \triangleq \partial_{\mathbf{f}}^2 \zeta^*. \quad (51)$$

595 \hat{Z}_{ε} and $\hat{\mathbf{x}}_{\text{a}}$ can be shown to exist; they can be read off from the explicit expressions of Z_{ε} and \mathbf{x}^{a} as functions of \mathbf{f} . These Jacobian and Hessians depend on the choice of the regularisation operator and they need to be computed analytically, which is simple but tedious, and not reported here since this is a regularisation-dependent calculation. The Hessian of the dual cost function Eq. (42) can then be written as the sum

$$\partial_{\mathbf{f}}^2 J_{\varepsilon}^* = \mathbf{\Lambda}_{\text{bo},\text{bo}} + \mathbf{H} \mathbf{\Omega}_{\text{bo},\text{bo}} \mathbf{H}^{\top}, \quad (52)$$

while the sensitivity matrix now reads

$$600 \quad \partial_{\mathbf{f}} \mathbf{x}^{\text{a}} = \mathbf{\Omega}_{\text{bo},\text{a}}^{\top} \mathbf{H}^{\top}. \quad (53)$$

Note that $\mathbf{\Omega}_{\text{bo},\text{bo}}$ can be interpreted as the covariance matrix of \mathbf{x} , the compounded observable vector as defined in Eq. (41) though seen as a random vector, on the assumption that \mathbf{x}^{b} and \mathbf{x}^{o} are connected via the W-barycentre \mathbf{x}^{a} and the optimal transference plans \mathbf{P}^{ba} and \mathbf{P}^{oa} , all seen as random vectors. Combining Eqs. (52,53) with Eq. (45), we finally obtain the expression for a factor \mathbf{X}^{a} of \mathbf{P}^{a} :

$$605 \quad \mathbf{X}^{\text{a}} = \mathbf{\Omega}_{\text{bo},\text{a}}^{\top} \mathbf{H}^{\top} \left[\mathbf{\Lambda}_{\text{bo},\text{bo}} + \mathbf{H} \mathbf{\Omega}_{\text{bo},\text{bo}} \mathbf{H}^{\top} \right]^{-1} \mathbf{\Lambda}_{\text{bo},\text{bo}}^{\frac{1}{2}}, \quad (54)$$

or, alternatively using the Sherman-Morrisson-Woodbury transformation, while assuming $\mathbf{\Omega}_{\text{bo},\text{bo}}$ to be invertible,

$$\mathbf{X}^{\text{a}} = \mathbf{\Omega}_{\text{bo},\text{a}}^{\top} \mathbf{\Omega}_{\text{bo},\text{bo}}^{-1} \left[\mathbf{\Omega}_{\text{bo},\text{bo}}^{-1} + \mathbf{H}^{\top} \mathbf{\Lambda}_{\text{bo},\text{bo}}^{-1} \mathbf{H} \right]^{-1} \mathbf{H}^{\top} \mathbf{\Lambda}_{\text{bo},\text{bo}}^{-\frac{1}{2}}. \quad (55)$$

610 These formulas are similar to the normal equations of classical DA. But mind that, in Eqs. (54,55), all the prior error statistics are encapsulated in $\mathbf{\Lambda}_{\text{bo},\text{bo}}$ whereas the impact of OT is encoded in $\mathbf{\Omega}_{\text{bo},\text{bo}}$. To be concrete, note that, when using Gaussian statistics Eq. (12), $\mathbf{\Lambda}_{\text{bo},\text{bo}}$ would simply read

$$\mathbf{\Lambda}_{\text{bo},\text{bo}} = \begin{bmatrix} \mathbf{\Lambda}_{\text{bb}} & \mathbf{0} \\ \mathbf{0} & \mathbf{\Lambda}_{\text{oo}} \end{bmatrix} = \begin{bmatrix} \mathbf{B} & \mathbf{0} \\ \mathbf{0} & \mathbf{R} \end{bmatrix}. \quad (56)$$

4.2 Interpretation

Further, we can perform a block decomposition of $\mathbf{\Omega}$ conformally to the spaces of \mathbf{x}^{b} and \mathbf{x}^{o} :

$$\mathbf{\Omega}_{\text{bo},\text{bo}} \triangleq \begin{bmatrix} \mathbf{\Omega}_{\text{bb}} & \mathbf{\Omega}_{\text{bo}} \\ \mathbf{\Omega}_{\text{bo}}^{\top} & \mathbf{\Omega}_{\text{oo}} \end{bmatrix}. \quad (57)$$

615 It can be shown that $\mathbf{\Omega}_{\text{bo}}$ is proportional to the optimal transference plan of the effective transport between \mathbf{x}^{o} and \mathbf{x}^{b} , and that the blocks of the diagonal are themselves diagonal and depend on the observable states:~

$$\mathbf{\Omega}_{\text{bo}} = \frac{1}{\varepsilon} \mathbf{P}^{\text{bo}}, \quad \mathbf{\Omega}_{\text{bb}} = \frac{1}{\varepsilon} \text{diag}(\mathbf{x}^{\text{b}}), \quad \mathbf{\Omega}_{\text{oo}} = \frac{1}{\varepsilon} \text{diag}(\mathbf{x}^{\text{o}}). \quad (58)$$

For instance, this could be shown by the explicit computation of $\Omega_{\text{bo},\text{bo}} = \varepsilon \partial_{\eta}^2 \hat{Z}_{\varepsilon}$.

Let us now examine the impact of OT on the analysis error covariance matrix. We first define

$$620 \quad \Delta = \Lambda_{\text{bo},\text{bo}}^{-\frac{1}{2}} \mathcal{H} \Omega_{\text{bo},\text{bo}}^{\frac{1}{2}}, \quad (59)$$

whose thin singular value decomposition is $\mathbf{U}\Sigma\mathbf{V}^{\top}$, where \mathbf{U} is an orthogonal matrix of size $(\mathfrak{N}_b + \mathfrak{N}_o) \times (\mathfrak{N}_b + \mathfrak{N}_o)$, Σ is a rectangular and diagonal matrix of size $(\mathfrak{N}_b + \mathfrak{N}_o) \times (N_b + N_o)$, and \mathbf{V} is an orthogonal matrix of size $(N_b + N_o) \times (N_b + N_o)$. Then, we *standardise* Eq. (54) following, e.g., Sect. 2.4.1 in Rodgers (2000):

$$\mathbf{X}^a = \Omega_{\text{bo},\text{a}}^{\top} \mathcal{H}^{\top} \left[\Lambda_{\text{bo},\text{bo}} + \mathcal{H} \Omega_{\text{bo},\text{bo}} \mathcal{H}^{\top} \right]^{-1} \Lambda_{\text{bo},\text{bo}}^{\frac{1}{2}}, \quad (60a)$$

$$625 \quad = \Omega_{\text{bo},\text{a}}^{\top} \mathcal{H}^{\top} \Lambda_{\text{bo},\text{bo}}^{-\frac{1}{2}} \left[\mathbf{I}_{\text{bo}} + \Lambda_{\text{bo},\text{bo}}^{-\frac{1}{2}} \mathcal{H} \Omega_{\text{bo},\text{bo}} \mathcal{H}^{\top} \Lambda_{\text{bo},\text{bo}}^{-\frac{1}{2}} \right]^{-1}, \quad (60b)$$

$$= \Omega_{\text{bo},\text{a}}^{\top} \Omega_{\text{bo},\text{bo}}^{-\frac{1}{2}} \Delta^{\top} \left[\mathbf{I}_{\text{bo}} + \Delta \Delta^{\top} \right]^{-1}, \quad (60c)$$

$$= \Omega_{\text{bo},\text{a}}^{\top} \Omega_{\text{bo},\text{bo}}^{-\frac{1}{2}} \mathbf{V}^{\top} \Sigma^{\top} \left[\mathbf{I}_{\text{bo}} + \Sigma \Sigma^{\top} \right]^{-1} \mathbf{U}^{\top}. \quad (60d)$$

Defining $\sigma = \left(\Sigma \Sigma^{\top} \right)^{\frac{1}{2}}$, which is square diagonal of size $(\mathfrak{N}_b + \mathfrak{N}_o) \times (\mathfrak{N}_b + \mathfrak{N}_o)$, we obtain, up to a multiplication by an irrelevant orthogonal matrix on the right, an equivalent factor for \mathbf{P}^a :

$$630 \quad \mathbf{X}^a = \Omega_{\text{bo},\text{a}}^{\top} \Omega_{\text{bo},\text{bo}}^{-\frac{1}{2}} \mathbf{V}^{\top} \frac{\sigma}{\mathbf{I}_{\text{bo}} + \sigma^2}. \quad (60e)$$

The diagonal values of σ , denoted $\sigma_i \geq 0$, represent the independent degrees freedom (dof) of information that can be extracted from the observations, which in our case is the first guess \mathbf{y}^b and the traditional observations \mathbf{y}^o , in contrast to Rodgers (2000) who only considers the dofs from \mathbf{y}^o . The higher the σ_i , the more information attached to the dof of index i , and the more squeezed the corresponding direction in \mathbf{X}^a and \mathbf{P}^a . From Eq. (59), and in particular its transpose: $\Delta^{\top} = \Omega_{\text{bo},\text{bo}}^{\frac{1}{2}} \mathcal{H}^{\top} \Lambda_{\text{bo},\text{bo}}^{-\frac{1}{2}}$

635 we can trace the flow of any piece of information. Such piece of information stems from the observation vectors, and hence its flow starts in Δ^{\top} from $\Lambda_{\text{bo},\text{bo}}^{-\frac{1}{2}}$ the square root of the precision matrix $\Lambda_{\text{bo},\text{bo}}^{-1}$. It is then transferred from the observation spaces to the observable spaces through \mathcal{H}^{\top} . It is finally optimally transported across the space of \mathbf{x}^b and \mathbf{x}^o by $\Omega_{\text{bo},\text{bo}}$ whose off-diagonal block is proportional to the transference plan \mathbf{P}^{bo} . Hence, OT is not a primary source of uncertainty, as \mathbf{y}^b and \mathbf{y}^o can be, but moves information in between the observable spaces.

640 Let us now check the OTDA analysis error covariance matrix \mathbf{P}^a in the classical DA limit. To that end, we study Eq. (54) in the classical limit. Similarly to Ω_{bb} and Ω_{oo} in Eq. (58), Ω_{aa} is defined as the covariance matrix of \mathbf{x}^a when only accounting for both OTs, and it can be shown that it reads

$$\Omega_{\text{aa}} = \frac{1}{\varepsilon} \text{diag}(\mathbf{x}^a). \quad (61)$$

When the cost tends to $\mathbf{C}_{\text{bo}}^{\infty}$, following the same arguments as in Sect. 2.4.3, \mathbf{x}^b , \mathbf{x}^o , and \mathbf{x}^a must merge and, consequently,

645 $\Omega_{\text{bo}} = \Omega_{\text{aa}} = \Omega_{\text{bb}} = \Omega_{\text{oo}}$. Hence, in this limit $\Omega_{\text{bo},\text{bo}} = \mathbf{1}_2 \Omega_{\text{aa}} \mathbf{1}_2^{\top}$, and $\Omega_{\text{bo},\text{a}} = \mathbf{1}_2 \Omega_{\text{aa}}$, with $\mathbf{1}_2 = [1 \quad 1]^{\top}$. Then, substitut-

ing these expressions of $\Omega_{\text{bo},\text{bo}}$ and $\Omega_{\text{bo},\text{a}}$ into Eq. (54), we get

$$\mathbf{X}^{\text{a}} = \Omega_{\text{aa}} \mathbf{1}_2^\top \mathcal{H}^\top \left[\Lambda_{\text{bo},\text{bo}} + \mathcal{H} \mathbf{1}_2 \Omega_{\text{aa}} \mathbf{1}_2^\top \mathcal{H}^\top \right]^{-1} \Lambda_{\text{bo},\text{bo}}^{\frac{1}{2}}, \quad (62\text{a})$$

$$= \Omega_{\text{aa}} \mathbf{1}_2^\top \mathcal{H}^\top \left[\mathbf{I}_{\text{bo}} + \Lambda_{\text{bo},\text{bo}}^{-1} \mathcal{H} \mathbf{1}_2 \Omega_{\text{aa}} \mathbf{1}_2^\top \mathcal{H}^\top \right]^{-1} \Lambda_{\text{bo},\text{bo}}^{-\frac{1}{2}}, \quad (62\text{b})$$

$$= \Omega_{\text{aa}} \left[\mathbf{I}_{\text{a}} + \mathbf{1}_2^\top \mathcal{H}^\top \Lambda_{\text{bo},\text{bo}}^{-1} \mathcal{H} \mathbf{1}_2 \Omega_{\text{aa}} \right]^{-1} \mathbf{1}_2^\top \mathcal{H}^\top \Lambda_{\text{bo},\text{bo}}^{-\frac{1}{2}}, \quad (62\text{c})$$

$$650 \quad = \left[\Omega_{\text{aa}}^{-1} + \mathbf{1}_2^\top \mathcal{H}^\top \Lambda_{\text{bo},\text{bo}}^{-1} \mathcal{H} \mathbf{1}_2 \right]^{-1} \mathbf{1}_2^\top \mathcal{H}^\top \Lambda_{\text{bo},\text{bo}}^{-\frac{1}{2}}, \quad (62\text{d})$$

where \mathbf{I}_{a} is the identity matrix of size N_{a} . From Eq. (62b) to Eq. (62c) we relied on the shift matrix lemma (e.g., Asch et al., 2016). For Ω_{aa}^{-1} in Eq. (62d) to exist, it must be assumed that $\mathbf{x}^{\text{a}} \in \mathcal{O}_{N_{\text{a}}}^{+,*}$, i.e. all the entries of \mathbf{x}^{a} are positive. This is verified when using entropic regularisation with $\varepsilon > 0$, no matter how small the entries of \mathbf{x}^{a} can be. Moreover, if \mathbf{x}^{a} has zero entries, \mathbf{x}^{a} can be represented as the limit of a sequence of positive discrete measures.

655 Now, since we have

$$\mathbf{A}^{-1} \triangleq \mathbf{1}_2^\top \mathcal{H}^\top \Lambda_{\text{bo},\text{bo}}^{-1} \mathcal{H} \mathbf{1}_2 = \Lambda_{\text{bb}}^{-1} + \mathbf{H}^\top \Lambda_{\text{oo}}^{-1} \mathbf{H}, \quad (63)$$

we conclude from Eq. (62d) that the classical limit of the analysis error covariance matrix is

$$\mathbf{P}^{\text{a}} = \mathbf{X}^{\text{a}} (\mathbf{X}^{\text{a}})^\top = \left[\Omega_{\text{aa}}^{-1} + \mathbf{A}^{-1} \right]^{-1} \mathbf{A}^{-1} \left[\Omega_{\text{aa}}^{-1} + \mathbf{A}^{-1} \right]^{-1}. \quad (64)$$

If the limit of \mathbf{x}^{a} when $\varepsilon \rightarrow 0^+$ is in $\mathcal{O}_{N_{\text{a}}}^{+,*}$, then $\Omega_{\text{aa}}^{-1} = \varepsilon \text{diag}(\mathbf{x}^{\text{a}})^{-1}$ must vanish. In this case:

$$660 \quad \mathbf{P}^{\text{a}} \xrightarrow{\varepsilon \rightarrow 0^+} \mathbf{A} = \left(\partial_{\mathbf{f}_b}^2 \zeta_b + \mathbf{H}^\top \partial_{\mathbf{f}_o}^2 \zeta_o \mathbf{H} \right)^{-1}, \quad (65)$$

which, assuming Gaussian errors, would read $\mathbf{P}^{\text{a}} = (\mathbf{B}^{-1} + \mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H})^{-1}$, as expected from classical DA. However, if some of the entries of \mathbf{x}^{a} vanish in the limit $\varepsilon \rightarrow 0^+$, we suspect that the limit of \mathbf{P}^{a} will be the classical analysis error covariance matrix \mathbf{A} but with the columns and rows associated to the vanishing entries of \mathbf{x}^{a} tapered to 0.

4.3 Numerical illustration

665 We consider the one-dimensional example where one half of the domain is observed, over $[0, \frac{1}{2}]$, $\mathbf{H} \in \mathcal{O}_{\mathfrak{N}_o \times N_o}^+$ with $\mathfrak{N}_o = N_o/2$ and $H_l^j = \delta_{l,j}$ for $l \in \llbracket 1, \mathfrak{N}_o \rrbracket$ and $j \in \llbracket 1, N_o \rrbracket$; the observations are unbalanced, $m(\mathbf{y}^{\text{b}}) = 1$ and $m(\mathbf{y}^{\text{o}}) = 1.49$; they have been generated through \mathcal{H} by discrete measures of mass 1 and 1.5, respectively; The cost metric has a quadratic dependence with the distance between sites, i.e. $[\mathbf{C}_{\text{ba}}]_{ik} = \lambda |r_i^{\text{b}} - r_k^{\text{a}}|^2$ and $[\mathbf{C}_{\text{oa}}]_{jk} = \lambda |r_j^{\text{o}} - r_k^{\text{a}}|^2$, where $\lambda = 10^3$. We use the results of Sect. 4.1 to compute the analysis error covariance matrix \mathbf{P}^{a} , the transference plan \mathbf{P}^{bo} , and the Jacobian $\Omega_{\text{bo},\text{a}}$. The numerical

670 results are displayed in Fig. 14. The OTDA analysis state is bimodal, some mass being left over to the right of the domain to account for the long tail of the first guess, which is far from the observation support. Hence, there is a vanishing field region, roughly $[0.6, 0.7]$, which separates the two components of the analysis state. As expected from OT theory, \mathbf{P}^{bo} seems to converge towards a (non-trivial and barely differentiable) Monge map which, in this discrete context, has two branches,

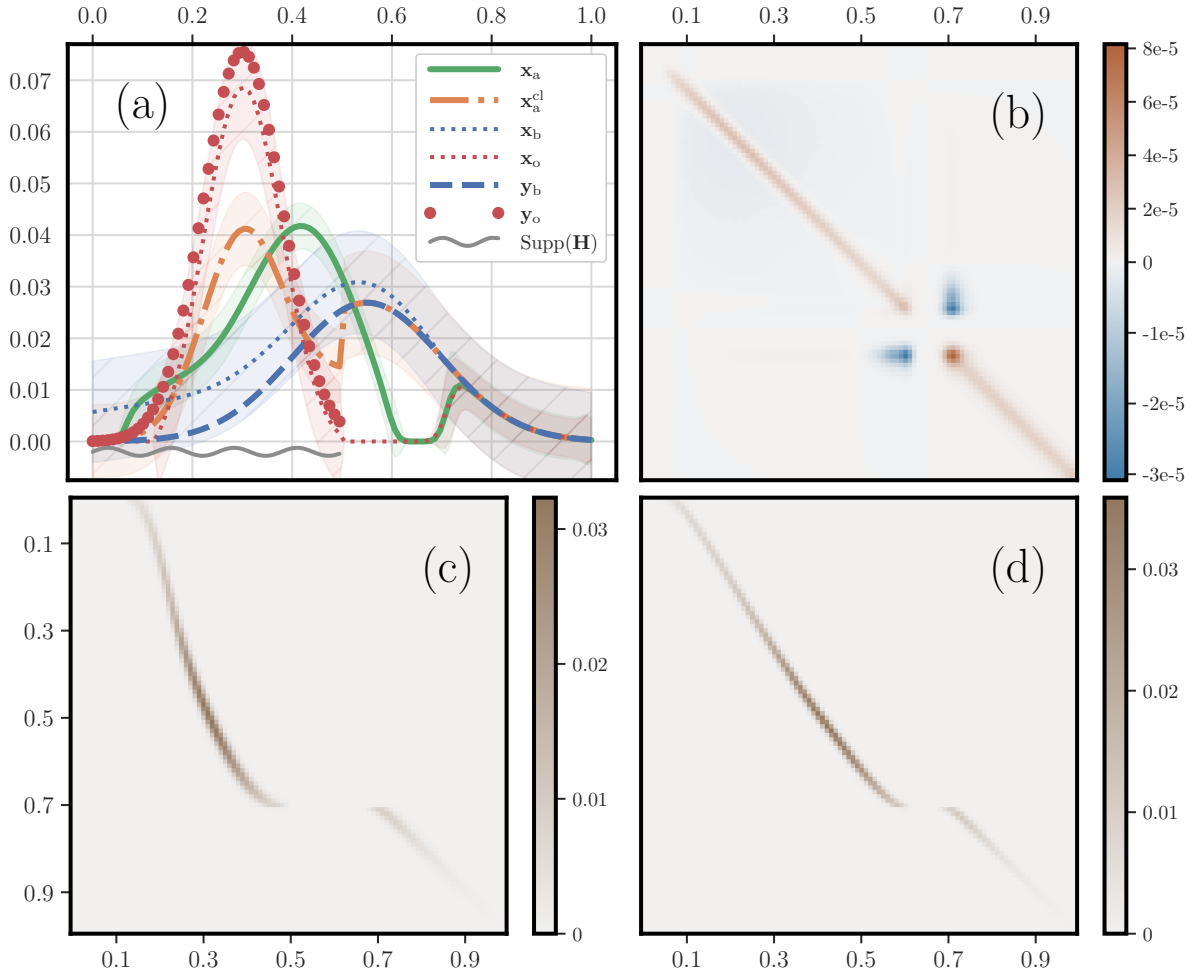


Figure 14. Illustration of the second-order analysis of an OTDA $\mathcal{3D-Var}$. Panel (a) shows the same plot as panel (a) of Fig. 11, but with errors bars (plus/minus the standard deviations) computed from the diagonal of the diagnosed posterior error covariance matrices associated to \mathbf{x}^a , \mathbf{x}^b , and \mathbf{x}^o . Panel (b) displays the analysis error covariance matrix \mathbf{P}^a . Panel (c) shows the optimal transference plan \mathbf{P}^{b^o} . Panel (d) shows the a, b block part of the Jacobian matrix $\Omega_{b^o, a}$, which is denoted Ω_{ba} .

separated by the gap created by the vanishing field region. The analysis error covariance matrix \mathbf{P}^a seems to converge to a
675 diagonal matrix, with the exception of the vanishing field region. Indeed, there seems to be an uncertainty as to how much mass
should be transferred from the first guess tail $[0.7, 1]$ to the main region $[0, 0.6]$. This is given away by peaks of variances near
the edges of the gap, and by negative covariances between the two edge points of the gap.

5 Conclusions

In this paper, we have introduced a theoretical framework for integrating nonlocal optimal transport (OT) metrics into data
680 assimilation (DA), which we refer to as *hybrid OTDA*. This framework addresses the inconsistencies initially identified by
Feyeux et al. when local metrics in classical DA are replaced with nonlocal ones based on OT.

Our focus has been on defining a ~~3D-Var~~3D-Var approach for hybrid OTDA and deriving the first- and second-order mo-
ments of its analysis. The hybrid OTDA ~~3D-Var~~3D-Var method blends classical DA and its background and observation error
685 statistics with a Wasserstein barycentre problem involving the observables associated with the first-guess and the observation
vector. Importantly, our work demonstrates that classical DA is encompassed within this theoretical framework.

We have shown that this optimisation problem can be decomposed and simplified into a hybrid OTDA problem with a single
OT problem based on an effective cost. This first problem yields the estimated \mathbf{x}^b and \mathbf{x}^o , followed by a pure W-barycentre
problem involving these states, whose solution is known as the McCann interpolant. This W-barycentre computation serves as
the final analysis step.

690 Our proposed method can be applied to sparsely and noisily observed systems, as expected from a robust DA method. It
can also accommodate non-trivial error statistics typical of a ~~3D-Var~~3D-Var approach. Furthermore, we have illustrated the
method's flexibility in defining cost metrics through various 1D and 2D numerical examples. We have empirically checked
how the OTDA analysis shifts towards the classical DA analysis, within the OTDA framework.

Note that, for now, some limitations apply; mainly the framework is presently meant for non-negative fields.

695 While we have looked into several other promising developments of our methodology, we have chosen not to report them in
this paper. These developments will be the subject of a future publication, including:

- the derivation of a Bayesian and probabilistic standpoint on OTDA,
- a generalised formalism where physical regularisation such as smoothness of the field can be enforced on the analysis
state,
- 700 – a *stochastic matrix* formalism, which is a substitute to using transference plans, but could offer more robustness in the
presence of entropic regularisation,
- employing cost matrices defined across several spaces, which is useful for realistic application where \mathbf{x}^b and \mathbf{x}^o lies
in very distinct spaces, such as the space of emission of a pollutant, and the space of the pollutant concentrations,
respectively.

705 While our primary focus in this paper was on the derivation and understanding of key cost functions within the hybrid OTDA
framework, we did not delve much into the numerical challenges, algorithmic complexity, or computing acceleration. For this
aspect of the developments, we would rather rely on the developments of the experts of OT who are continuously improving
on the efficiency of numerical OT (e.g., Flamary et al., 2021).

In addition to strengthening the developments mentioned above, our future research will explore the application of the hybrid
710 OTDA formalism in a sequential DA framework, as this paper concentrated solely on the static analysis. We are also interested

in investigating the role played by error statistics and cost metrics $\{\zeta_b, \zeta_o, \mathbf{C}_{ba}, \mathbf{C}_{oa}\}$ and their balancing in the hybrid OTDA analysis, as well as developing their objective tuning.

Code availability. The products of this paper are exclusively optimisation problems and methods to solve them; their implementation (code) used in the illustrative sections rely on freely available software to solve the optimisation problems, mainly L-BFGS-B and its implementation in Scipy <https://github.com/scipy/scipy> and the Python Optimal Transport library and its implementation <https://github.com/PythonOT>.

Appendix A: From the primal to the dual cost function for the full problem

This appendix is dedicated to the derivation of the transformation from a Lagrangian variant of the primal problem to the dual cost function Eq. (20). It takes the form of the following series of transformations of the problem, from a Lagrangian to a dual cost function:

$$\begin{aligned}
720 \quad \mathcal{L} = & \min_{\mathbf{x}^b \mathbf{x}^o \mathbf{x}^a} \left[\max_{\mathbf{f}_b} \{(\mathbf{y}^b - \mathbf{x}^b)^\top \mathbf{f}_b - \zeta_b^*(\mathbf{f}_b)\} + \max_{\mathbf{f}_o} \{(\mathbf{y}^o - \mathbf{H}\mathbf{x}^o)^\top \mathbf{f}_o - \zeta_o^*(\mathbf{f}_o)\} \right. \\
& + \min_{\mathbf{P} \in \mathcal{O}_{b,o,a}^+} \{P_{ijk} C_{ba}^{ik} + P_{ijk} C_{oa}^{jk}\} \\
& \left. + \max_{\mathbf{h}_b, \mathbf{h}_o, \mathbf{f}_a} \{h_i^b (x_i^b - P_{ijk} 1_o^j 1_a^k) + h_j^o (x_j^o - P_{ijk} 1_b^i 1_a^k) + f_a^k (x_k^a - P_{ijk} 1_b^i 1_o^j)\} \right], \tag{A1a}
\end{aligned}$$

$$\begin{aligned}
= & \min_{\mathbf{x}^a \mathbf{x}^b \mathbf{x}^o} \left[\max_{\mathbf{f}_b, \mathbf{f}_o} \{(\mathbf{y}^b - \mathbf{x}^b)^\top \mathbf{f}_b - \zeta_b^*(\mathbf{f}_b) + (\mathbf{y}^o - \mathbf{H}\mathbf{x}^o)^\top \mathbf{f}_o - \zeta_o^*(\mathbf{f}_o)\} \right. \\
& + \max_{\mathbf{h}_b, \mathbf{h}_o, \mathbf{f}_a} \min_{\mathbf{P} \in \mathcal{O}_{b,o,a}^+} \{P_{ijk} C_{ba}^{ik} + P_{ijk} C_{oa}^{jk}\} \\
725 \quad & \left. + h_i^b (x_i^b - P_{ijk} 1_o^j 1_a^k) + h_j^o (x_j^o - P_{ijk} 1_b^i 1_a^k) + f_a^k (x_k^a - P_{ijk} 1_b^i 1_o^j)\} \right], \tag{A1b}
\end{aligned}$$

$$\begin{aligned}
= & \max_{\mathbf{h}_b, \mathbf{h}_o, \mathbf{f}_a} \min_{\mathbf{f}_b, \mathbf{f}_o} \min_{\mathbf{x}^a \mathbf{x}^b \mathbf{x}^o} \left[(\mathbf{y}^b - \mathbf{x}^b)^\top \mathbf{f}_b - \zeta_b^*(\mathbf{f}_b) + (\mathbf{y}^o - \mathbf{H}\mathbf{x}^o)^\top \mathbf{f}_o - \zeta_o^*(\mathbf{f}_o) \right. \\
& \left. + P_{ijk} C_{ba}^{ik} + P_{ijk} C_{oa}^{jk} + h_i^b (x_i^b - P_{ijk} 1_o^j 1_a^k) + h_j^o (x_j^o - P_{ijk} 1_b^i 1_a^k) + f_a^k (x_k^a - P_{ijk} 1_b^i 1_o^j) \right], \tag{A1c}
\end{aligned}$$

$$\begin{aligned}
= & \max_{\mathbf{f}_b, \mathbf{f}_o} \{ \mathbf{f}_b^\top \mathbf{y}^b + \mathbf{f}_o^\top \mathbf{y}^o - \zeta_b^*(\mathbf{f}_b) - \zeta_o^*(\mathbf{f}_o) \} \\
& + \min_{\mathbf{h}_b, \mathbf{h}_o, \mathbf{f}_a} \left[(\mathbf{h}_b - \mathbf{f}_b)^\top \mathbf{x}^b + (\mathbf{h}_o - \mathbf{H}^\top \mathbf{f}_o)^\top \mathbf{x}^o + \mathbf{f}_a^\top \mathbf{x}^a \right. \\
730 \quad & \left. + P_{ijk} (C_{ba}^{ik} + C_{oa}^{jk} - h_i^b 1_o^j 1_a^k - h_j^o 1_b^i 1_a^k - f_a^k 1_b^i 1_o^j) \right], \tag{A1d}
\end{aligned}$$

$$\begin{aligned}
= & \max_{\mathbf{f}_b, \mathbf{f}_o} \left[\mathbf{f}_b^\top \mathbf{y}^b + \mathbf{f}_o^\top \mathbf{y}^o - \zeta_b^*(\mathbf{f}_b) - \zeta_o^*(\mathbf{f}_o) \right. \\
& \left. + \min_{\mathbf{P} \in \mathcal{O}_{b,o,a}^+} P_{ijk} \left(C_{ba}^{ik} + C_{oa}^{jk} - f_b^i 1_o^j 1_a^k - H_l^j f_o^l 1_b^i 1_a^k \right) \right]. \tag{A1e}
\end{aligned}$$

In Eq. (A1a), the maps ζ_b^* and ζ_o^* are the Legendre-Fenchel transforms of the maps ζ_b and ζ_o , respectively. From Eq. (A1d) to Eq. (A1e), taking the minimum over the observables \mathbf{x}^b , \mathbf{x}^o , and \mathbf{x}^a implies to enforce $\mathbf{h}_b = \mathbf{f}_b$, $\mathbf{h}_o = \mathbf{H}^\top \mathbf{f}_o$, and $\mathbf{f}_a = \mathbf{0}$.

735 Hence, we obtain the dual problem which only depends on the Lagrange multipliers:

$$\mathcal{L}^* = \max_{(\mathbf{f}_b, \mathbf{f}_o) \in \mathcal{U}_{b_o}^*(\mathbf{C}_{ba}, \mathbf{C}_{oa}, \mathbf{H})} \{ \mathbf{f}_b^\top \mathbf{y}^b + \mathbf{f}_o^\top \mathbf{y}^o - \zeta_b^*(\mathbf{f}_b) - \zeta_o^*(\mathbf{f}_o) \}, \quad (\text{A2a})$$

where the $*$ symbol refers to *dual* and where the polyhedron $\mathcal{U}_{b_o}^*(\mathbf{C}_{ba}, \mathbf{C}_{oa}, \mathbf{H})$ is defined by

$$\mathcal{U}_{b_o}^*(\mathbf{C}_{ba}, \mathbf{C}_{oa}, \mathbf{H}) \triangleq \left\{ \mathbf{f}_b \in \mathbb{R}^{\mathfrak{N}_b}, \mathbf{f}_o \in \mathbb{R}^{\mathfrak{N}_o} : \forall i, j, k, \quad f_b^i + f_o^l H_l^j \leq C_{ba}^{ik} + C_{oa}^{jk} \right\}. \quad (\text{A2b})$$

740 The inequality constraints of the polyhedron $\mathcal{U}_{b_o}^*$ stem from the positivity constraint $P_{ijk} \geq 0$ in Eq. (A1e). Very importantly, we have the coincidence of the minimum of the primal problem with the maximum of the dual problem $\mathcal{L} = \mathcal{L}^*$, a property called *strong duality* (see Sect. 5.2 in Boyd and Vandenberghe, 2004). Strong duality can for instance be achieved if both the primal and dual cost functions are convex, which is the case here.

Appendix B: Derivation of the two-step hybrid optimal transport data assimilation algorithm

Here we derive the two-step algorithm elaborated in Sec. 2.4.2. Moreover, entropic regularisation is added to the problem.

745 B1 First step: simplified hybrid optimal transport data assimilation problem

The first step of the full OTDA algorithm is a simplified OTDA problem based on a single OT problem driven by the cost \mathbf{C}_{b_o} . The corresponding high-level primal cost function is

$$\mathcal{L} = \min_{\mathbf{x}^b \in \mathcal{O}_b^+, \mathbf{x}^o \in \mathcal{O}_o^+} \{ \zeta_b(\mathbf{y}^b - \mathbf{x}^b) + \zeta_o(\mathbf{y}^o - \mathbf{H}\mathbf{x}^o) + W_{\mathbf{C}_{b_o}}(\mathbf{x}^b, \mathbf{x}^o) \}. \quad (\text{B1})$$

The associated (lower level) primal cost function, but adding entropic regularisation ($\varepsilon > 0$), is then

$$750 \quad \mathcal{L}_\varepsilon = \min_{\mathbf{x}^b \in \mathcal{O}_b^+, \mathbf{x}^o \in \mathcal{O}_o^+} \left[\zeta_b(\mathbf{y}^b - \mathbf{x}^b) + \zeta_o(\mathbf{y}^o - \mathbf{H}\mathbf{x}^o) + \min_{\mathbf{P} \in \mathcal{U}_{b_o}} \left(\varepsilon \mathcal{K}(\mathbf{P} | \boldsymbol{\nu}) + P_{ij} C_{b_o}^{ij} \right) \right]. \quad (\text{B2a})$$

In this optimisation problem, the admissible set of transference plans, i.e. the set of all 2-tensors of non negative entries whose marginals are \mathbf{x}^b and \mathbf{x}^o , is defined by

$$\mathcal{U}_{b_o} \triangleq \left\{ \mathbf{P} \in \mathcal{O}_{b_o}^+ : \mathbf{P} \mathbf{1}_o = \mathbf{x}^b, \quad \mathbf{P}^\top \mathbf{1}_b = \mathbf{x}^o \right\}. \quad (\text{B2b})$$

755 Since \mathbf{x}^b and \mathbf{x}^o are not predetermined, the prior transference plan $\boldsymbol{\nu}$ cannot be selected from \mathcal{U}_{b_o} a priori. The simplest choice, which we decided to implement, is hence to set ν_{ij} to a constant, which assumes some statistical prior independence of \mathbf{x}^b and \mathbf{x}^o . A derivation of the dual problem equivalent to \mathcal{L}_ε can be obtained in the exact same way as in the previous subsection, although it is now less cluttered since there is only one OT to account for, instead of two. The associated Lagrangian is

$$\begin{aligned} \mathcal{L}_\varepsilon = & \max_{\mathbf{f}_b \in \mathbb{R}^{\mathfrak{N}_b}, \mathbf{f}_o \in \mathbb{R}^{\mathfrak{N}_o}} \left[\mathbf{f}_b^\top \mathbf{y}^b + \mathbf{f}_o^\top \mathbf{y}^o - \zeta_b^*(\mathbf{f}_b) - \zeta_o^*(\mathbf{f}_o) \right. \\ & \left. + \min_{\mathbf{P} \in \mathcal{O}_{b_o}^+} \left(\varepsilon \sum_{ij} \left\{ P_{ij} \ln \frac{P_{ij}}{\nu_{ij}} - P_{ij} + \nu_{ij} \right\} + P_{ij} \left\{ C_{b_o}^{ij} - f_b^i \nu_{ij} - H_l^j f_o^l \nu_{ij} \right\} \right) \right]. \quad (\text{B3}) \end{aligned}$$

760 Again, the maps ζ_b^* and ζ_o^* are the Legendre-Fenchel transforms of the maps ζ_b and ζ_o . The variables \mathbf{f}_b and \mathbf{f}_o are Lagrange vectors; they are used to enforce the marginals of the transference plan associated to $W_{C_{b,o}}$. The unconstrained minimisation over \mathbf{P} , i.e. the inner minimisation problem in Eq. (B3), is obtained by cancelling the gradient with respect to \mathbf{P} , which yields

$$P_{ij} = \nu_{ij} e^{(f_b^i + f_o^j H_i^j - C_{b,o}^{ij})/\varepsilon}. \quad (\text{B4})$$

Substituting this solution into minus the Lagrangian $-\mathcal{L}_\varepsilon$ gives the regularised dual problem

$$\mathcal{J}_\varepsilon^* = \min_{\mathbf{f}_b \in \mathbb{R}^{n_b} \mathbf{f}_o \in \mathbb{R}^{n_o}} J_\varepsilon^*(\mathbf{f}_b, \mathbf{f}_o), \quad (\text{B5a})$$

765 with the associated Lagrangian

$$J_\varepsilon^*(\mathbf{f}_b, \mathbf{f}_o) = \varepsilon (Z_\varepsilon - \mathbf{m}(\boldsymbol{\nu})) + \zeta_b^*(\mathbf{f}_b) + \zeta_o^*(\mathbf{f}_o) - \mathbf{f}_b^\top \mathbf{y}^b - \mathbf{f}_o^\top \mathbf{y}^o, \quad (\text{B5b})$$

which relies on the *partition function*

$$Z_\varepsilon = \sum_{ij} P_{ij}. \quad (\text{B5c})$$

770 The notation $\mathcal{J}_\varepsilon^*$ and J_ε^* , rather than $\mathcal{L}_\varepsilon^*$ and L_ε^* , signifies that we work on the opposite of $\mathcal{L}_\varepsilon^*$ and L_ε^* so as to obtain a dual problem to be minimised rather than maximised. Most importantly, we have, under conditions that will be satisfied in the following, the coincidence of the two minima $\mathcal{J}_\varepsilon^* = -\mathcal{L}_\varepsilon$, i.e. strong duality. Assuming one can obtain a proper correspondence between the optimal $\mathbf{f}_b, \mathbf{f}_o$ of the dual problem and $\mathbf{x}^b, \mathbf{x}^o$ of the primal problem, this implies, once again, that the primal problem can be traded for the dual problem.

775 Even though the regularised optimisation problem is slightly different from the unregularised one, a difference which is controlled by the value of ε , the new dual optimisation problem is free, i.e. without constraints. It can be solved as it is, using for instance the L-BFGS-B minimiser (Liu and Nocedal, 1989). The advantage of the regularised dual formulation is two-fold: the dual cost function is unconstrained (free optimisation) and we will trade a minimisation over $N_b \times N_o$ variables for a minimisation over $N_b + N_o$ variables. This dual formulation can be viewed as a generalised *Physical-space Statistical Analysis System* (PSAS) formalism (Courtier, 1997), an approach where classical DA algebra is mostly carried out in observation space.

780

Once the optimal values for \mathbf{f}_b and \mathbf{f}_o are obtained, the optimal discrete Kantorovich transference plan \mathbf{P} can be computed using Eq. (B4). As a result, as marginals of this transference plan, the solutions for the observables are

$$x_i^b = P_{ij} 1_o^j = \sum_j P_{ij}, \quad x_j^o = P_{ij} 1_b^i = \sum_i P_{ij}. \quad (\text{B6})$$

B2 Second step: Wasserstein barycentre

785 Now that we have obtained the observables \mathbf{x}^b and \mathbf{x}^o via Eq. (B6), we would like to compute their W-barycentre. The joint mass m of these observables can be computed:

$$m = \mathbf{m}(\mathbf{x}^b) = \mathbf{m}(\mathbf{x}^o). \quad (\text{B7})$$

The high-level primal cost function of this W-barycentre problem is

$$J_w = \min_{\mathbf{x}^a \in \mathcal{O}_{N_a}^+} \{W_{C_{ba}}(\mathbf{x}^b, \mathbf{x}^a) + W_{C_{oa}}(\mathbf{x}^o, \mathbf{x}^a)\}. \quad (\text{B8})$$

790 We have found and practised several ways to solve this problem. One way is to compute the McCann interpolant. This is theoretically elegant but Eq. (26) did not leverage regularisation of the W-barycentre problem. Instead, the approach reported here is to use the dual optimisation problem, in conjunction with the entropic regularisation at finite $\varepsilon > 0$. We leverage our knowledge of the mass m resulting from the first step of the algorithm by enforcing the mass in the cost function, $m(\mathbf{P}) = m$. This seems redundant but it actually yields *by construction* and very naturally a numerical efficient algorithm comparable to
 795 the *ad hoc* log-domain scheme proposed in Sect. 4.4 of Peyré and Cuturi (2019).

Again, one way, though not the only one, to write the primal problem goes through the use of a gluing transference plan, a 3-tensor whose marginals are \mathbf{x}^b , \mathbf{x}^o , and \mathbf{x}^a :

$$L_\varepsilon = \min_{\mathbf{x}^a \in \mathcal{O}_{N_a}^+, \mathbf{P} \in \mathcal{U}_{\text{boa}}(\mathbf{x}^a)} \{ \mathbf{P} \cdot \mathbf{C}_{\text{boa}} + \varepsilon \mathcal{K}(\mathbf{P} | \nu) + \mathbf{f}_b^\top \mathbf{x}^b + \mathbf{f}_o^\top \mathbf{x}^o \}, \quad (\text{B9a})$$

where $[\mathbf{C}_{\text{boa}}]_{ijk} = C_{ba}^{ik} + C_{oa}^{jk}$, the binary operator \cdot denotes the contraction of tensors, and

$$800 \quad \mathcal{U}_{\text{boa}}(\mathbf{x}^a) \triangleq \left\{ \mathbf{P} \in \mathcal{O}_{b,o,a}^+ : \forall i, P_{ijk} 1_o^j 1_a^k = x_i^b, \quad \forall j, P_{ijk} 1_b^i 1_a^k = x_j^o, \quad \forall k, P_{ijk} 1_b^i 1_o^j = x_k^a \right\}. \quad (\text{B9b})$$

The 3-tensor ν is chosen to be $\nu_{ijk} = x_i^b x_j^o / (m N_a)$, which is uniform in k and for which $m(\nu) = m$. The resulting dual problem is

$$\mathcal{J}^* = \min_{\mathbf{f}_b \in \mathbb{R}^{n_b}, \mathbf{f}_o \in \mathbb{R}^{n_o}} J^*(\mathbf{f}_b, \mathbf{f}_o), \quad (\text{B10a})$$

where the associated Lagrangian is

$$805 \quad J^*(\mathbf{f}_b, \mathbf{f}_o) = \varepsilon \left(m \ln \frac{Z_\varepsilon}{m} + m - m(\nu) \right) - \mathbf{f}_b^\top \mathbf{x}^b - \mathbf{f}_o^\top \mathbf{x}^o, \quad (\text{B10b})$$

with the partition function

$$Z_\varepsilon = \sum_{ijk} \nu_{ijk} e^{(f_b^i + f_o^j H_l^j - C_{ba}^{ik} - C_{oa}^{jk}) / \varepsilon}. \quad (\text{B10c})$$

This partition function is elegant but impractical since in high dimension a 3-tensor might be too large to store and compute with. However the partition function Eq. (B10c) can be simplified by noticing that

$$810 \quad Z_\varepsilon = \sum_{ij} \nu_{ij} e^{(f_b^i + f_o^j H_l^j - C_{bo}^{ij}) / \varepsilon}, \quad (\text{B11})$$

where we introduced the effective cost metric

$$[\mathbf{C}_{bo}]_{ij} = -\varepsilon \ln \left(\sum_k \frac{\nu_{ijk}}{\nu_{ij}} e^{-(C_{ba}^{ik} + C_{oa}^{jk}) / \varepsilon} \right), \quad (\text{B12})$$

815 which is the regularised cost – known in statistics and machine learning as a *soft-plus* transform – of Eq. (22c). The 2-tensor ν_{ij} plays the same role as that of the first step of the algorithm; we choose it as $\nu_{ij} = x_i^b x_j^o / m$, for which $m(\boldsymbol{\nu}) = m$. The dual problem now only involves 2-tensors and becomes numerically more efficient. Given the optimal \mathbf{f}_b and \mathbf{f}_o , the (glued) optimal transference plan \mathbf{P}^{boa} is formally given by

$$P_{ijk}^{\text{boa}} = \frac{\nu_{ijk}}{Z_\varepsilon} e^{(f_b^i + f_o^l H_l^j - C_{ba}^{ik} - C_{oa}^{jk})/\varepsilon}. \quad (\text{B13})$$

The W-barycentre \mathbf{x}^a is then given as a marginal of \mathbf{P}^{boa} :

$$x_k^a = P_{ijk}^{\text{boa}} 1_b^i 1_o^j = \frac{1}{Z_\varepsilon} \sum_{ij} \nu_{ijk} e^{(f_b^i + f_o^l H_l^j - C_{ba}^{ik} - C_{oa}^{jk})/\varepsilon}. \quad (\text{B14})$$

820 Because of the normalisation of the transference plan to m , the entropic regularisation exhibits a $\varepsilon m \ln Z_\varepsilon$ instead of $\varepsilon Z_\varepsilon$. This systematically enforces normalisation in the computations of the gradients, as well as in the course of the numerical optimisation of the dual cost function, *de facto* working in log-domain. We experienced more stable computations and the ability to reach smaller ε , as compared to the case without normalisation. This completes the solution through the 2-step OTDA algorithm.

825 *Author contributions.* MB, PJV and AF developed the methodology. MB implemented the numerics. MB wrote the manuscript. MB, PJV, AF, JDLB and YR revised the manuscript.

Competing interests. The contact author has declared that none of the authors has any competing interests.

Acknowledgements. The authors would like to thank two anonymous Reviewers for their remarks and suggestions that helped shape the paper, as well as the Executive and Handling Editor O. Talagrand. This research has been supported by the national research project ANR-830 ARGONAUT (grant no. ANR-19-CE01-0007, PollutAnt and gReenhouse Gases emissiOns moNitoring from spAcE at high ResolUTion). CERA is a member of Institut Pierre-Simon Laplace (IPSL).

References

- Amodei, M. and Stein, J.: Deterministic and fuzzy verification methods for a hierarchy of numerical models, *Meteorological Applications*, 16, 191–203, <https://doi.org/10.1002/met.101>, 2009.
- 835 Asch, M., Bocquet, M., and Nodet, M.: *Data Assimilation: Methods, Algorithms, and Applications*, *Fundamentals of Algorithms*, SIAM, Philadelphia, ISBN 978-1-611974-53-9, <https://doi.org/10.1137/1.9781611974546>, 2016.
- Bocquet, M.: Towards optimal choices of control space representation for geophysical data assimilation, *Mon. Wea. Rev.*, 137, 2331–2348, <https://doi.org/10.1175/2009MWR2789.1>, 2009.
- Bocquet, M., Wu, L., and Chevallier, F.: Bayesian design of control space for optimal assimilation of observations. I: Consistent multiscale
840 formalism, *Q. J. R. Meteorol. Soc.*, 137, 1340–1356, <https://doi.org/10.1002/qj.837>, 2011.
- Boyd, S. P. and Vandenberghe, L.: *Convex optimization*, Cambridge university press, ISBN 978-0521833783, 2004.
- Briggs, W. M. and Levine, R. A.: Wavelets and field forecast verification, *Mon. Wea. Rev.*, 125, 1329–1341, [https://doi.org/1520-0493\(1997\)125<1329:WAFFV>2.0.CO;2](https://doi.org/1520-0493(1997)125<1329:WAFFV>2.0.CO;2), 1997.
- Carrassi, A., Bocquet, M., Bertino, L., and Evensen, G.: Data Assimilation in the Geosciences: An overview on methods, issues, and perspectives, *WIREs Climate Change*, 9, e535, <https://doi.org/10.1002/wcc.535>, 2018.
845
- Chizat, L., Peyré, G., Schmitzer, B., and Vialard, F.-X.: Scaling algorithms for unbalanced optimal transport problems, *Mathematics of Computation*, 87, 2563–2609, <https://doi.org/10.1090/mcom/3303>, 2018.
- Courtier, P.: Dual formulation of four-dimensional variational assimilation, *Q. J. R. Meteorol. Soc.*, 123, 2449–2461, <https://doi.org/10.1002/qj.49712354414>, 1997.
- 850 Daley, R.: *Atmospheric Data Analysis*, Cambridge University Press, New-York, 1991.
- Davis, C., Brown, B., and Bullock, R.: Object-based verification of precipitation forecasts. Part I: Methodology and application to mesoscale rain areas, *Mon. Wea. Rev.*, 134, 772–1784, <https://doi.org/10.1175/MWR3146.1>, 2006a.
- Davis, C., Brown, B., and Bullock, R.: Object-based verification of precipitation forecasts. Part II: Application to convective rain systems, *Mon. Wea. Rev.*, 134, 1785–1795, <https://doi.org/10.1175/MWR3145.1>, 2006b.
- 855 Duc, L. and Sawada, Y.: Geometry of rainfall ensemble means: from arithmetic averages to Gaussian-Hellinger barycenters in unbalanced optimal transport, *J. Meteor. Soc. Japan*, 102, 35–47, <https://doi.org/10.2151/jmsj.2024-003>, 2024.
- Ebert, E. E.: Fuzzy verification of high-resolution gridded forecasts: a review and proposed framework, *Meteorological applications*, 15, 51–64, <https://doi.org/10.1002/met.25>, 2008.
- El Moselhy, T. A. and Marzouk, Y. M.: Bayesian inference with optimal maps, *J. Comp. Phys.*, 231, 7815–7850,
860 <https://doi.org/10.1016/j.jcp.2012.07.022>, 2012.
- Evensen, G., Vossepoel, F. C., and van Leeuwen. P. J.: *Data Assimilation Fundamentals: A Unified Formulation of the State and Parameter Estimation Problem*, *Springer Textbooks in Earth Sciences, Geography and Environment*, Springer Cham, ISBN 978-3-030-96708-6, <https://doi.org/10.1007/978-3-030-96709-3>, 2022.
- Farchi, A. and Bocquet, M.: Review article: Comparison of local particle filters and new implementations, *Nonlin. Processes Geophys.*, 25,
865 765–807, <https://doi.org/10.5194/npg-25-765-2018>, 2018.
- Farchi, A., Bocquet, M., Rouston, Y., Mathieu, A., and Quérel, A.: Using the Wasserstein distance to compare fields of pollutants: application to the radionuclide atmospheric dispersion of the Fukushima-Daiichi accident, *Tellus B*, 68, 31682, <https://doi.org/10.3402/tellusb.v68.31682>, 2016.

- Feyeux, N.: Transport optimal pour l'assimilation de données images, Ph.D. thesis, Université Grenoble Alpes, <https://inria.hal.science/tel-01480695>, 2016.
- 870
- Feyeux, N., Vidard, A., and Nodet, M.: Optimal transport for variational data assimilation, *Nonlin. Processes Geophys.*, 25, 55–66, <https://doi.org/10.5194/npg-25-55-2018>, 2018.
- Flamary, R., Courty, N., Gramfort, A., Alaya, M. Z., Boisbunon, A., Chambon, S., Chapel, L., Corenflos, A., K., F., Fournier, N., Gautheron, L., Gayraud, N. T. H., Janati, H., Rakotomamonjy, A., Redko, I., Rolet, A., Schutz, A., Seguy, V., Sutherland, D. J., Tavenard, R., Tong, A.,
- 875 and Vayer, T.: POT: Python Optimal Transport, *Journal of Machine Learning Research*, 22, 1–8, <http://jmlr.org/papers/v22/20-451.html>, 2021.
- Gangbo, W. and McCann, R. J.: The geometry of optimal transportation, *Acta Mathematica*, 177, 113–1618, <https://doi.org/10.1007/BF02392620>, 1996.
- Gilleland, E., Ahijevych, D. A., Brown, B. G., and Ebert, E. E.: Verifying forecasts spatially, *Bull. Amer. Meteor. Soc.*, 91, 1365–1373,
- 880 <https://doi.org/10.1175/2010BAMS2819.1>, 2010a.
- Gilleland, E., Lindström, J., and Lindgren, F.: Analyzing the image warp forecast verification method on precipitation fields from the ICP, *Weather and Forecasting*, 25, 1249–1262, <https://doi.org/10.1175/2010WAF2222365.1>, 2010b.
- Hoffman, R. N. and Grassotti, C.: A Technique for Assimilating SSM/I Observations of Marine Atmospheric Storms: Tests with ECMWF Analyses, *Journal of Applied Meteorology and Climatology*, 35, 1177–1188, [https://doi.org/10.1175/1520-](https://doi.org/10.1175/1520-0450(1996)035<1177:ATFASO>2.0.CO;2)
- 885 [0450\(1996\)035<1177:ATFASO>2.0.CO;2](https://doi.org/10.1175/1520-0450(1996)035<1177:ATFASO>2.0.CO;2), 1996.
- Hoffman, R. N., Liu, Z., Louis, J.-F., and Grassotti, C.: Distortion representation of forecast errors, *Mon. Wea. Rev.*, 123, 2758–2770, [https://doi.org/10.1175/1520-0493\(1995\)123<2758:DROFE>2.0.CO;2](https://doi.org/10.1175/1520-0493(1995)123<2758:DROFE>2.0.CO;2), 1995.
- Janjić, T., Bormann, N., Bocquet, M., Carton, J. A., Cohn, S. E., Dance, S. L., Losa, S. N., Nichols, N. K., Potthast, R., Waller, J. A., and Weston, P.: On the representation error in data assimilation, *Q. J. R. Meteorol. Soc.*, 144, 1257–1278, <https://doi.org/10.1002/qj.3130>,
- 890 2018.
- Kalnay, E.: *Atmospheric Modeling, Data Assimilation and Predictability*, Cambridge University Press, Cambridge, 2003.
- Keil, C. and Craig, G. C.: A displacement and amplitude score employing an optical flow technique, *Weather and Forecasting*, 24, 1297–1308, <https://doi.org/10.1175/2009WAF2222247.1>, 2009.
- Lack, S. A., Limpert, G. L., and Fox, N. I.: An object-oriented multiscale verification scheme, *Weather and Forecasting*, 25, 79–92,
- 895 <https://doi.org/10.1175/2009WAF2222245.1>, 2010.
- Le Coz, C., Tantet, A., Flamary, R., and Plougonven, R.: Optimal transport for the multi-model combination of sub-seasonal ensemble forecasts, *EGU23-13445*, <https://doi.org/10.5194/egusphere-egu23-13445>, 2023.
- Liu, D. C. and Nocedal, J.: On the limited memory BFGS method for large scale optimization, *Mathematical Programming*, 45, 503–528, <https://doi.org/10.1007/BF01589116>, 1989.
- 900 Lledó, L., Skok, G., and Haiden, T.: Estimating location errors in precipitation forecasts with the Wasserstein and Attribution distances, *EMS2023-602*, <https://doi.org/10.5194/ems2023-602>, 2023.
- Marzouk, Y., Moselhy, T., Parno, M., and Spantini, A.: An introduction to sampling via measure transport, in: *Handbook of Uncertainty Quantification*, edited by Ghanem, R., Higdon, D., and Owhadi, H., chap. 23, pp. 785–825, Springer International Publishing, Cham, https://doi.org/10.1007/978-3-319-12385-1_23, 2017.
- 905 Monge, G.: Mémoire sur la théorie des déblais et des remblais, in: *Histoire de l'Académie Royale des Sciences de Paris*, pp. 666–704, 1781.

- Necker, T., Wolfgruber, L., Kugler, L., Weissmann, M., Dorninger, M., and Serafin, S.: The fractions skill score for ensemble forecast verification, <https://doi.org/10.22541/au.169169008.89657659/v1>, 2023.
- Ning, L. and Carli, F. P., Ebtehaj, A. M., Fofoula-Georgiou, E., and Georgiou, T. T.: Coping with model error in variational data assimilation using optimal mass transport, *Water Resources Research*, 50, 5817–5830, <https://doi.org/10.1002/2013WR014966>, 2014.
- 910 Oliver, D. S.: Minimization for conditional simulation: Relationship to optimal transport, *J. Comp. Phys.*, 265, 1–15, <https://doi.org/10.1016/j.jcp.2014.01.048>, 2014.
- Peyré, G. and Cuturi, M.: Computational Optimal Transport: With Applications to Data Science, *Foundations and Trends® in Machine Learning*, 11, 355–607, <https://doi.org/10.1561/22000000073>, 2019.
- Plu, M.: A variational formulation for translation and assimilation of coherent structures, *Nonlin. Processes Geophys.*, 20, 793–801,
 915 <https://doi.org/10.5194/npg-20-793-2013>, 2013.
- Ravela, S., Emanuel, K., and McLaughlin, D.: Data assimilation by field alignment, *Physica D*, 230, 127–145, <https://doi.org/10.1016/j.physd.2006.09.035>, 2007.
- Rodgers, C. D.: Inverse methods for atmospheric sounding, vol. 2, World Scientific, Series on Atmospheric, Oceanic and Planetary Physics, ISBN 978-981-02-2740-1, <https://doi.org/10.1142/3171>, 2000.
- 920 Sinkhorn, R.: A relationship between arbitrary positive matrices and doubly stochastic matrices, *The annals of mathematical statistics*, 35, 876–879, <http://www.jstor.org/stable/2238545>, 1964.
- Skok, G.: Precipitation attribution distance, *Atmospheric Research*, 295, 106998, <https://doi.org/10.1016/j.atmosres.2023.106998>, 2023.
- Talagrand, O.: Assimilation of Observations, an Introduction, *J. Meteor. Soc. Japan*, 75, 191–209, https://doi.org/10.2151/jmsj1965.75.1B_191, 1997.
- 925 Tamang, S. K., Ebtehaj, A., Zou, D., and Lerman, G.: Regularized variational data assimilation for bias treatment using the Wasserstein metric, *Q. J. R. Meteorol. Soc.*, 146, 2332–2346, <https://doi.org/10.1002/qj.3794>, 2020.
- Tamang, S. K., Ebtehaj, A., van Leeuwen, P. J., Zou, D., and Lerman, G.: Ensemble Riemannian data assimilation over the Wasserstein space, *Nonlin. Processes Geophys.*, 28, 295–309, <https://doi.org/10.5194/npg-28-295-2021>, 2021.
- Tamang, S. K., Ebtehaj, A., van Leeuwen, P. J., Lerman, G., and Fofoula-Georgiou, E.: Ensemble Riemannian Data Assimilation: Towards
 930 High-dimensional Implementation, *Nonlin. Processes Geophys.*, 29, 77–92, <https://doi.org/10.5194/npg-29-77-2022>, 2022.
- Vanderbecken, P. J., Dumont Le Brazidec, J., Farchi, A., Bocquet, M., Roustan, Y., Potier, E., and Broquet, G.: Accounting for meteorological biases in simulated plumes using smarter metrics, *Atmos. Meas. Tech.*, 16, 1745–1766, <https://doi.org/10.5194/amt-16-1745-2023>, 2023.
- Vilani, C.: Topics in Optimal Transportation, vol. 58 of *Graduate Studies in Mathematics*, American Mathematical Society, Providence, Rhode Island, 2003.
- 935 Vilani, C.: Optimal Transport: Old and New, vol. 338 of *Die Grundlehren der Mathematischen Wissenschaften*, Springer-Verlag, Berlin Heidelberg, 2009.
- Ying, Y.: A Multiscale Alignment Method for Ensemble Filtering with Displacement Errors, *Mon. Wea. Rev.*, 147, 4553–4565, <https://doi.org/10.1175/MWR-D-19-0170.1>, 2019.
- Ying, Y., Anderson, J. L., and Bertino, L.: Improving Vortex Position Accuracy with a New Multiscale Alignment Ensemble Filter, *Mon.
 940 Wea. Rev.*, 151, 1387–405, <https://doi.org/10.1175/MWR-D-22-0140.1>, 2023.
- Zhou, W., Bovik, A. C., Sheikh, H. R., and Simoncelli, E.: Image quality assessment: from error visibility to structural similarity, *IEEE Transactions on Image Processing*, 13, 600–612, <https://doi.org/10.1109/TIP.2003.819861>, 2004.