



HAL
open science

Observation data compression for variational assimilation of dynamical systems

Sibo Cheng, Didier Lucor, Jean-Philippe Argaud

► **To cite this version:**

Sibo Cheng, Didier Lucor, Jean-Philippe Argaud. Observation data compression for variational assimilation of dynamical systems. *Journal of computational science*, 2021, 53, pp.101405. 10.1016/j.jocs.2021.101405 . hal-03335014

HAL Id: hal-03335014

<https://hal.science/hal-03335014>

Submitted on 2 Aug 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

Observation data compression for variational assimilation of dynamical systems

Sibo Cheng^{1,2,3}, Didier Lucor², Jean-Philippe Argaud³

¹ *Data Science Institute, Department of Computing, Imperial College London, UK*

² *Université Paris-Saclay, CNRS, Laboratoire Interdisciplinaire des Sciences du Numérique, France*

³ *EDF R&D Saclay, France*

Abstract

Accurate estimation of error covariances (both background and observation) is crucial for efficient observation compression approaches in data assimilation of large-scale dynamical problems. We propose a new combination of a covariance tuning algorithm with existing PCA-type data compression approaches, either observation- or information-based, with the aim of reducing the computational cost of real-time updating at each assimilation step. Relying on a local assumption of flow-independent error covariances, dynamical assimilation residuals are used to adjust the covariance in each assimilation window. The estimated covariances then contribute to better specify the principal components of either the observation dynamics or the state-observation sensitivity. The proposed approaches are first validated on a shallow water twin experiment with correlated and non-homogeneous observation error. Proper selection of flow-independent assimilation windows, together with sampling density for background error estimation, and sensitivity of the approaches to the observations error covariance knowledge, are also discussed and illustrated with various numerical tests and results. The method is then applied to a more challenging industrial hydrological model with real-world data and non-linear transformation operator provided by an operational precipitation-flow simulation software.

Keywords: Data assimilation, Observation compression, Error covariance estimation, Information entropy, Hydrological application

1. Introduction

Data assimilation (DA) is applied in a wide range of industrial problems, such as numerical weather prediction (NWP) [1], hydrology, fire forecasting [2] or nuclear engineering [3]. Recently, DA methods have also been used to COVID-19 pandemic analysis, including predicting disease diffusion and proposing optimal vaccination strategies ([4]). DA algorithms are often used in dynamical systems for continuously updating state estimation/prediction. They have recently made their way to other fields such as biomedical applications [5] or quantitative economics [6]. These methods rely on a weighted combination of different sources of noisy information, including prior numerical estimation (also known as background states) and real-time observations, to improve field reconstruction or parameters calibration. DA methods are often used to deal with problems of large dimensions, especially in NWP [7], [8] (up to 10^9) or in geoscience [9], leading to computational difficulty for real-time updating, if not infeasible. Several strategies for optimizing the computational cost have been developed, including graph-based domain localization [10], observations selection [11], matrix decomposition [12] or reduced-order Kalman Filter [13]. It is also a common practice to combine DA algorithms with classical dynamical system reduction techniques, such as the Proper Orthogonal Decomposition (POD) or the Empirical Interpolation Method (EIM e.g. [14]). Most of these methods rely on either precise knowledge of state variables (e.g. modes in POD) or strong prior assumptions (e.g. cut-off radius in domain localization [15]). Meanwhile, with the increase of available observation precision in DA applications, the observation data compression via low-rank approximation methods

21 has been continuously studied for alleviating the computational cost, especially in a sequential
22 data assimilation chain. These methods, which consist of extracting principal information in
23 observation data, have been widely applied in various branches of engineering, especially for
24 high dimensional problems. An important advantage of observation compression, regarding
25 other methods that directly reduce the state space dimension, is that no extra operation/
26 knowledge of the state dynamics is required, making the compression error more controllable
27 and estimable. Two classical compression methods are discussed and implemented in this work:
28 the POD-type projection by extracting principal components in the observation dynamic [16]
29 and the information-based compression based on the information entropy analysis [8]. The
30 latter aims to select the most impacting observations to the analyzed state by calculating the
31 prior-posterior information entropy gap. Since the noises are introduced by prior errors in DA
32 systems, the information entropy estimation relies on both background and observation error
33 covariance matrices.

34
35 For both observation- and information-based approaches, the data compression is carried out
36 with a noise-normalized dataset [7], [8]. The knowledge of prior error covariances thus becomes
37 crucial for applying these methods. However, the specification of these covariances, especially
38 the background matrix, remains one of the most challenging problems in data assimilation due
39 to the high dimension of the problem and limited prior data [17],[18]. Much attention was
40 given to improving the error covariance specification in dynamical data assimilation models,
41 particularly by the meteorological society. Several methods have been developed to this end,
42 such as the NMC approach [19], the DI01 [20] iterative method and the Desroziers estimation
43 [21]. In this paper, we focus on the latest. Unlike some other methods (e.g. [20], [18]), the
44 Desroziers estimation does not depend on the specific structure of the error covariances, and
45 it provides a non-parametric estimation of full covariances as output of the algorithm. Based
46 on the residual analysis in variational assimilation, this approach has been widely applied in
47 industrial problems, especially in NWP. Recent works of [22] prove its convergence in the ideal
48 case. Another considerable strength of the Desroziers estimation is that dynamic residual data
49 can be used for the covariance estimation. For this reason, a huge ensemble size is not required
50 for high dimensional problems, unlike, for instance, in the NMC method.

51
52 In this paper, based on the Desroziers estimation, we have introduced the concept of piece-
53 wise covariance estimation for both observation- and information-based compression strategies.
54 We apply the Desroziers method to estimate error covariances in a fixed time range, also known
55 as the flow-independent window where the error covariances are supposed to be time-invariant.
56 Therefore, the choice of the flow-independent window and the residual samplings play an es-
57 sential role in this algorithm. The window size should be sufficiently long to gather enough
58 time-variant sampling but not too long to consider the error covariances, especially the back-
59 ground matrix, being constant.

60
61 The observation- and information-based (with piecewise covariances estimation) data com-
62 pression are first implemented in a twin experiment framework using 2D shallow water equations
63 with a linear transformation operator. The observation covariance is supposed to be perfectly
64 known *a priori*. The two approaches with different choices of flow-independent windows are
65 compared in this model while changing the truncation parameter. Numerical results show that
66 the observation-based (POD-type) compression is in general over-performed by the information-
67 based approach and that a non-balanced sampling in piecewise covariance estimation results in
68 a less optimal compression. We then apply these methods to a real-world hydrological model
69 to improve river flow prediction/reanalysis by correcting historical daily precipitation measures
70 [23]. Both the precipitation and the river flow data are spatially distributed. The physical
71 simulation is performed using the operating MORDOR-TS software [24], developed by EDF
72 and the study area is around the Tarn river, in the south of France. The precipitation-flow

73 simulation is carried out through conceptual watersheds modeling, which ensures its high com-
74 putational efficiency. In this hydrological application, both the background and the observation
75 matrices are estimated using the Desroziers method with daily observed flow data for around
76 10 years (1990 to 2000). Results show that in this industrial application where both \mathbf{B} and \mathbf{R}
77 are not well known, the performance of the information-based strategy is similar to the one of
78 observation-based.

79

80 The paper is organized as follows. In section 2, the principle and the notation of data
81 assimilation are briefly introduced. We then introduce the observation- and information-based
82 compression strategies in section 3. The applications of 2D shallow water twin experiments and
83 an industrial hydrological model are shown respectively in section 5 and 6. We finish the paper
84 with a discussion.

85 2. Variational data assimilation

The objective of data assimilation algorithms is to improve the estimation of some physical fields or parameters \mathbf{x} based on two sources of information: a prior simulation/forecast \mathbf{x}_b and an observation vector \mathbf{y} . The theoretical value of the current state is denoted by a vector \mathbf{x}_{true} , also known as the true state. Variational DA algorithms aim to find an optimally weighted compromise between the prior estimation \mathbf{x}_b and the observation \mathbf{y} by minimising the cost function J defined as

$$J(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}_b)^T \mathbf{B}^{-1}(\mathbf{x} - \mathbf{x}_b) + \frac{1}{2}(\mathbf{y} - \mathcal{H}(\mathbf{x}))^T \mathbf{R}^{-1}(\mathbf{y} - \mathcal{H}(\mathbf{x})) \quad (1)$$

$$= \frac{1}{2}\|\mathbf{x} - \mathbf{x}_b\|_{\mathbf{B}^{-1}}^2 + \frac{1}{2}\|\mathbf{y} - \mathcal{H}(\mathbf{x})\|_{\mathbf{R}^{-1}}^2 \quad (2)$$

where \mathcal{H} denotes the transformation operator from the state space to one of the observations. \mathbf{B} and \mathbf{R} are the associated error covariance matrices, i.e.

$$\mathbf{B} = \text{Cov}(\epsilon_b, \epsilon_b), \quad \mathbf{R} = \text{Cov}(\epsilon_y, \epsilon_y), \quad (3)$$

where

$$\epsilon_b = \mathbf{x}_b - \mathbf{x}_{\text{true}}, \quad \epsilon_y = \mathcal{H}(\mathbf{x}_{\text{true}}) - \mathbf{y}. \quad (4)$$

Thus the inverse of these covariance matrices (i.e. $\mathbf{B}^{-1}, \mathbf{R}^{-1}$) represents the weights of these two information sources in the objective function. Prior errors ϵ_b, ϵ_y are supposed to be centered Gaussian, characterised by the error covariance matrices, i.e.

$$\epsilon_b \sim \mathcal{N}(0, \mathbf{B}), \quad \epsilon_y \sim \mathcal{N}(0, \mathbf{R}). \quad (5)$$

The optimization problem of Eq. 1, so called three-dimensional variational (3D-Var) formulation, is a general representation of variational assimilation while the model error is not considered. The output of Eq. 1 is denoted as \mathbf{x}_a , i.e.

$$\mathbf{x}_a = \underset{\mathbf{x}}{\text{argmin}} \left(J(\mathbf{x}) \right). \quad (6)$$

If \mathcal{H} can be approximated by some linear operator \mathbf{H} , Eq. 6 can be solved via BLUE (Best Linearized Unbiased Estimator) formulation,

$$\mathbf{x}_a = \mathbf{x}_b + \mathbf{K}(\mathbf{y} - \mathbf{H}\mathbf{x}_b) \quad (7)$$

$$\mathbf{A} = (\mathbf{I} - \mathbf{K}\mathbf{H})\mathbf{B} \quad (8)$$

where $\mathbf{A} = \text{Cov}(\mathbf{x}_a - \mathbf{x}_{\text{true}})$ is the analyzed error covariance and the \mathbf{K} matrix, given by

$$\mathbf{K} = \mathbf{B}\mathbf{H}^T(\mathbf{H}\mathbf{B}\mathbf{H}^T + \mathbf{R})^{-1} \quad (9)$$

86 is so called the Kalman gain matrix. In the rest of this paper, we denote \mathbf{H} as the linearized
 87 transformation operator. The case when \mathcal{H} is non-linear is more challenging for finding the
 88 minimum of Eq. 1, especially for high-dimensional problems. The resolution involves often gra-
 89 dient descent algorithms (relying on algorithms such as "L-BFGS-B" [25] and on adjoint-based
 90 [11] numerical techniques.

91

Variational assimilation algorithms could be applied to dynamical systems through sequen-
 tial applications using a transition operator $\mathcal{M}_{t^k \rightarrow t^{k+1}}$ (from time t^k to t^{k+1}), where

$$\mathbf{x}_{t^{k+1}} = \mathcal{M}_{t^k \rightarrow t^{k+1}}(\mathbf{x}_{t^k}). \quad (10)$$

The forecasting thus depends on the knowledge of transition operator $\mathcal{M}_{t^k \rightarrow t^{k+1}}$ and the cor-
 rected state at the current time \mathbf{x}_{a,t^k} . Typically, the current background state is often given by
 the forecasting from the previous step, i.e.

$$\mathbf{x}_{b,t^k} = \mathcal{M}_{t^{k-1} \rightarrow t^k}(\mathbf{x}_{a,t^{k-1}}). \quad (11)$$

92 Obviously, a more accurate reanalysis $\mathbf{x}_{a,t^{k-1}}$ leads to a more reliable forecasting \mathbf{x}_{b,t^k} . It
 93 is known that as long as the transformation operator \mathcal{H} and the transition operator \mathcal{M} are
 94 linear, **the analysis based on the variational method and the Kalman filter results in the same**
 95 **forecasting [9], for dynamical (4D-Var) assimilation problems.** Theoretically, the evolution of
 96 the \mathbf{B} matrix could also be estimated thanks to the transition operator. However, in practice,
 97 the perfect knowledge of \mathcal{M} is often unavailable. Much attention is given to quantify the model
 98 error in assimilation, for example, in weak-constraint 4D-VAR [26]. Recent work of [27] involves
 99 deep learning techniques to improve the estimation of $\mathcal{M}_{t^{k-1} \rightarrow t^k}$.

100 3. Observation data compression

101 DA algorithms are often used to perform real-time corrections of dynamical systems with
 102 large dimensions, leading to an essential requirement of computational efficiency. In this work,
 103 we are interested in a low-rank approximations of the observation vector which can reduce the
 104 cost of real-time updating in DA algorithms.

105 3.1. Observation-based compression (OC)

The works of [28] and [29] are based on a PCA-type reduction of the observation dynamics.
 More precisely, a set of n_{obs} observation snapshots is represented by a matrix $\mathbf{Y} \in \mathbb{R}_{[\dim(\mathbf{y}) \times n_{obs}]}$
 where each column $\mathbf{Y}[:, i]$ represents an individual observation vector of dimension m at a fixed
 time t_i , i.e.

$$\mathbf{Y}[:, i] = \mathbf{y}_{t=t_i}. \quad (12)$$

106 Thus \mathbf{Y} describes the evolution of the observation vector \mathbf{y} including observation error. We
 107 work with the error-normalized data $\mathbf{R}^{-1/2}\mathbf{Y}$ [7] whose empirical covariance \mathbf{C} can be written
 108 and decomposed as

$$\mathbf{C} = \frac{1}{n_{obs} - 1} \mathbf{R}^{-1/2} \mathbf{Y} \mathbf{Y}^T \mathbf{R}^{-1/2} = \tilde{\mathbf{L}} \tilde{\mathbf{D}} \tilde{\mathbf{L}}^T \quad (13)$$

where the columns of $\tilde{\mathbf{L}}$ are the principal components and $\tilde{\mathbf{D}}$ represents the associated eigen-
 values in a decreasing order. This decomposition is known as the principal component analysis
 (PCA) decomposition. We can construct a projection operator $\tilde{\mathbf{L}}_q$ with minimum loss of in-
 formation (represented by eigenvalues in the covariance matrix) by simply keeping the q first
 columns in $\tilde{\mathbf{L}}$. q is also known as the truncation parameter. In fact, this projection operator

can also be obtained by a singular value decomposition (SVD), without computing the full covariance matrix \mathbf{C} , i.e.

$$\mathbf{R}^{-1/2}\mathbf{Y} = \tilde{\mathbf{L}}_q \tilde{\Sigma} \tilde{\mathbf{V}}_q^T \quad (14)$$

where $\tilde{\mathbf{L}}_q$ and $\tilde{\mathbf{V}}_q$ are orthogonal matrices, i.e. $\tilde{\mathbf{L}}_q^T \tilde{\mathbf{L}}_q = \tilde{\mathbf{V}}_q^T \tilde{\mathbf{V}}_q = \mathbf{I}$ and $\tilde{\Sigma} \tilde{\Sigma}^T = \tilde{\mathbf{D}}$ since all eigenvalues are non negative. The assumption is made for the observation error covariances to be constant (flow-independent), which is a common practice in data assimilation (e.g [7]). For each DA optimization, instead of updating with the full observation vector \mathbf{y} , the correction is made with the reduced observation

$$\tilde{\mathbf{y}}_q = \tilde{\mathbf{L}}_q^T \mathbf{R}^{-1/2} \mathbf{y}. \quad (15)$$

The new observation error covariance $\tilde{\mathbf{R}}$ and the new state-observation transformation operator $\tilde{\mathcal{H}}$ can be written as

$$\tilde{\mathbf{R}}_q = \tilde{\mathbf{L}}_q^T \mathbf{R}^{-1/2} \mathbf{R} \mathbf{R}^{-1/2} \tilde{\mathbf{L}}_q = \mathbf{I}_q, \quad \tilde{\mathcal{H}}_q = \tilde{\mathbf{L}}_q^T \mathbf{R}^{-1/2} \circ \mathcal{H}. \quad (16)$$

109 The DA algorithm can then be performed on $(\mathbf{x}_b, \tilde{\mathbf{y}}_q, \mathbf{B}, \tilde{\mathbf{R}}_q, \tilde{\mathcal{H}}_q)$ instead of $(\mathbf{x}_b, \mathbf{y}, \mathbf{B}, \mathbf{R}, \mathcal{H})$.
 110 This method could be seen as a classical POD approach applied to error-normalised observation
 111 data by extracting modes of higher variances against time. It is pointed out by [28] and [7]
 112 that performing PCA on noise-normalised observation data can improve the method efficiency
 113 and reduce the impact of observation error during the compression procedure.

114 3.2. Information-based compression (IC)

The observation-based data reduction retains the principal directions of the observation dynamic. However, these directions are not necessarily the most impacting in state correction. A continuous effort has been devoted to quantify and compute the sensitivity of the analysis states to the observations (e.g. [11]), which leads to a more refined observation compression in DA. More precisely, this sensitivity may be expressed by the influence matrix \mathbf{S} [30], defined as

$$\mathbf{S} = \frac{\partial \mathcal{H}(\mathbf{x}_a)}{\partial \mathbf{x}_a} = \mathbf{K}^T \mathbf{H}^T. \quad (17)$$

According to [8], the information given by the influence matrix can be roughly quantified via two indicators, the degree of freedom for signal (DFS) which represents the prior-posterior mutual information and the entropy reduction (ER) which represents the evolution of Shannon information content, respectively defined as

$$\text{DFS} = \mathbb{E}[(\mathbf{x}_a - \mathbf{x}_b)^T \mathbf{B}^{-1} (\mathbf{x}_a - \mathbf{x}_b)] = \text{Tr}(\mathbf{S}) \quad (18)$$

$$\text{ER} = H(\mathbf{x}) - H(\mathbf{x}|\mathbf{y}) = -\frac{1}{2} \ln \left(\det(\mathbf{I} - \mathbf{S}) \right) \quad (19)$$

where H is the entropy of a distribution, noted here $H(x)$ for simplicity. Eqs. 18 and 19 are derived for a centred Gaussian vector \mathbf{x} . For both measures, we observe that observations associated with the largest eigenvalues of \mathbf{S} have the greatest information content. Using an intermediate matrix $\mathbf{M} = \mathbf{R}^{-1/2} \mathbf{H} \mathbf{B}^{1/2}$, Eqs. 18-19 could be rewritten as

$$\text{DFS} = \text{Tr}(\mathbf{M} \mathbf{M}^T (\mathbf{I} + \mathbf{M} \mathbf{M}^T)^{-1}) \quad (20)$$

$$\text{ER} = \frac{1}{2} \ln \left(\det(\mathbf{I} + \mathbf{M} \mathbf{M}^T) \right). \quad (21)$$

As stated in the work of [31], the observation projection operator which minimizes the information loss is given by $\hat{\mathbf{L}}_q \mathbf{R}^{-1/2}$, where $\hat{\mathbf{L}}_q$ is the matrix whose columns contain the eigenvectors of $\mathbf{M} \mathbf{M}^T = \mathbf{R}^{-1/2} \mathbf{H} \mathbf{B} \mathbf{H}^T \mathbf{R}^{-1/2}$. DA algorithms could then be performed with

$$\hat{\mathbf{y}}_q = \hat{\mathbf{L}}_q^T \mathbf{R}^{-1/2} \mathbf{y}, \quad \hat{\mathbf{R}}_q = \mathbf{I}_q, \quad \hat{\mathcal{H}}_q = \hat{\mathbf{L}}_q^T \mathbf{R}^{-1/2} \circ \mathcal{H}. \quad (22)$$

115 We remind that, from the computational point of view, the only difference between the OC and
 116 IC is the way the low-rank projection \mathbf{L}_q is obtained. For both approaches, the specification of
 117 error covariance matrices (either background or observation) is crucial to provide an efficient
 118 compression. On the other hand, data compression strategies can reduce the computational
 119 cost of covariance tuning methods, especially for multidimensional and multivariate problems.
 120 Therefore, the precise knowledge of \mathbf{HBH}^T and \mathbf{R} is crucial for this method. **However, as**
 121 **pointed out by [8], the condition number of the analysis covariance matrix \mathbf{A} can be higher when**
 122 **using IC approach compared to performing DA with the full observation data set. Therefore,**
 123 **the risk of matrix ill-conditioning is worth monitoring when applying this compression method.**
 124

125 3.3. Optimal truncation parameter for compression methods

The determination of the truncated parameter q , i.e. number of modes kept in the reduced space, is crucial in data compression. The choice of the threshold often depends on available data [32]. Several criteria were considered, such as the information losing rate E_q and the matrix conditioning *a posteriori* μ_q , defined as

$$E_q = \frac{\|\Sigma - \phi_q\|_\infty}{\|\Sigma\|_\infty} = 1 - \frac{\sigma_{q-1}}{\sigma_q} \quad (23)$$

$$\mu_q = \frac{\sigma_1}{\sigma_q} \quad (24)$$

where Σ is the diagonal matrices with all eigenvalues of the covariance matrix and $\sigma_{i,i=1..}$ represent the associated real eigenvalues in the decreasing order of absolute value. According to the study of [33], an optimal choice of the truncation parameter can be obtained by combining the two previous indicators, with an objective function f , defined as

$$f(\sigma_q) = E_q + \mu_q = \frac{\sigma_q \sigma_{q-1} + \sigma_1}{\sigma_1 \sigma_q}. \quad (25)$$

126 Assuming $\|\sigma_q - \sigma_1\| \gg \|\sigma_q - \sigma_{q-1}\|$, one could easily prove that Eq. 25 achieves the minimum
 127 when $\sigma_q = \sqrt{\sigma_1}$. With this choice, we manage to both reduce the matrix ill-conditioning and
 128 remove less significant modes, as proved in real-world DA application[33]. Another advantage
 129 of this criteria is that the computation of the full spectrum of covariances is not required. By
 130 applying Lanczos-type methods [34], we can stop the algorithm when the current eigenvalue is
 131 inferior to $\sqrt{\sigma_1}$.

132 4. Piecewise estimation of error covariances

The \mathbf{R} matrix is required for both OC and IC approaches. Furthermore, the construction of $\hat{\mathbf{L}}_q^T$ in IC requires a precise knowledge of the matrix production \mathbf{HBH}^T . However, the knowledge of both matrices often remains challenging in data assimilation [17]. Continuous effort was devoted to improve the error covariances specification [35], [18]. A classical approach based on residual analysis, and later a more complete version are respectively given by [36] and [21]. They show that under the assumption of flow-independent error covariance, i.e. \mathbf{B} and \mathbf{R} being invariant against time in a certain period, the following equations hold

$$\mathbf{R} = \mathbb{E} \left[\left(\mathbf{y} - \mathcal{H}(\mathbf{x}_a) \right) \left(\mathbf{y} - \mathcal{H}(\mathbf{x}_b) \right)^T \right] \quad (26)$$

$$\mathbf{HBH}^T = \mathbb{E} \left[\left(\mathbf{y} - \mathcal{H}(\mathbf{x}_b) \right) \left(\mathbf{y} - \mathcal{H}(\mathbf{x}_b) \right)^T \right] - \mathbf{R}. \quad (27)$$

Under these hypothesis, combining Eq. 26 and 27 leads to

$$\mathbf{HBH}^T = \mathbb{E} \left[\left(\mathcal{H}(\mathbf{x}_a) - \mathcal{H}(\mathbf{x}_b) \right) \left(\mathbf{y} - \mathcal{H}(\mathbf{x}_b) \right)^T \right]. \quad (28)$$

133 In order to somewhat alleviate these strong hypotheses, a simple idea is to take the expectation
 134 operators in Eq. 26, 27 and 28 in assimilation windows where the flow-independent assumption
 135 stands, resulting in a piecewise estimation of both \mathbf{B} and \mathbf{R} . More precisely, a sequence of
 136 estimated background matrices \mathbf{B}_{T_i} could be computed via residual covariances, where T_i refer
 137 to flow-independent periods of \mathbf{B} in a dynamical system. In other words, \mathbf{B} is considered as
 138 invariant between $t = T_i$ and $t = T_{i+1}$. The estimation of \mathbf{R}_{T_i} , if required, follows the same
 139 principle using Eq. 26. When the knowledge of \mathbf{R} matrix is precise *a priori*, the estimation of
 140 Eq. 27 is privileged because of its lower computational cost since no evaluation of the analyzed
 141 state \mathbf{x}_a is required. According to [36], when the observation error is dominated by background
 142 error (i.e. $\text{Tr}(\mathbf{R}) \ll \text{Tr}(\mathbf{B})$), \mathbf{HBH}^T can be estimated directly by $\mathbb{E}\left[\left(\mathbf{y} - \mathcal{H}(\mathbf{x}_b)\right)\left(\mathbf{y} - \mathcal{H}(\mathbf{x}_b)\right)^T\right]$.

By definition,

$$\mathbf{HBH}^T = \mathbb{E}\left[\left(\mathcal{H}(\mathbf{x}_{\text{true}}) - \mathcal{H}(\mathbf{x}_b)\right)\left(\mathcal{H}(\mathbf{x}_{\text{true}}) - \mathcal{H}(\mathbf{x}_b)\right)^T\right] \quad (29)$$

143 represents the background error covariances projected in the observation space. Therefore the
 144 information-based observation compression which is based on a PCA-type analysis, can also
 145 be interpreted as a projection of \mathbf{y} along the directions where the background errors are most
 146 important. Recently, it was also reported in the literature (e.g. [37]) that the convergence
 147 (towards the exact observation matrix) of the iterative method can still be ensured when the
 148 background and observation error correlation length-scales are similar, which was contrary
 149 to what was previously thought [38]. Although this innovation-based covariance estimation
 150 approach has been widely applied in DA applications, some drawbacks have also been noticed.
 151 For example, the application of this method in real problems often requires post-processing
 152 of the \mathbf{R} matrix. It is shown in [23] that the regularized matrix may converge to some other
 153 solution rather than the exact observation matrix.

154 5. Shallow water twin experiments

155 5.1. Experiments set up

156 For evaluating the performance of different data compression approaches, we set up a twin
 157 experiment framework with a simplified 2D shallow water dynamical model which is frequently
 158 used for testing data assimilation algorithms (e.g [11], [18]). A cylinder of water is positioned
 159 in the middle of the study field of size $20\text{mm} \times 20\text{mm}$ and released at the initial time $t = 0\text{s}$
 160 (i.e. with no initial speed), leading to a non-linear wave-propagation. The dynamics of the
 161 water level h (in mm), as well as horizontal and vertical velocity (in 0.1m/s) field (respectively
 162 denoted as u and v), is given by the non-conservative shallow water equations

$$\begin{aligned} \frac{\partial u}{\partial t} &= -g \frac{\partial}{\partial x}(h) - bu \\ \frac{\partial v}{\partial t} &= -g \frac{\partial}{\partial y}(h) - bv \\ \frac{\partial h}{\partial t} &= -\frac{\partial}{\partial x}(uh) - \frac{\partial}{\partial y}(vh) \\ u_{t=0} &= 0 \\ v_{t=0} &= 0 \end{aligned} \quad (30)$$

where $b = 0.1$ is the viscous drag coefficient and the earth gravity constant g is thus scaled to 1. These equations are discretized in a 20×20 regular grid, solved by first-order finite difference method with a time discretization $\delta_t = 10^{-4}\text{s}$. This resolution is considered as the reference (i.e. the true state \mathbf{x}_{true}) latter when performing DA algorithms. The state variables in this DA

modeling are the combination of the velocity fields $\{u\}_{20 \times 20}$ and $\{v\}_{20 \times 20}$. The evolution of the reference ($\mathbf{x}_{\text{true},t}$) state is illustrated in Fig. 1. Spatially correlated prior error is then generated artificially for simulating the background state with a standard deviation $\sigma_{b,0} = 0.2$, i.e.

$$\mathbf{x}_{b,t=0} \sim \mathcal{N}(\mathbf{x}_{\text{true},t=0}, \mathbf{B}_{t=0}) \quad \text{where} \quad \mathbf{B}_{t=0} = \sigma_{b,0}^2 \text{corr}(\mathbf{B}). \quad (31)$$

The background error correlation matrix $\text{corr}(\mathbf{B})$ is set to be isotropic (rotational invariant), following the second-order auto-aggressive (SOAR, also known as Balgovind) function,

$$\phi_{\mathbf{B}}(r) = \left(1 + \frac{r}{L_{\mathbf{B}}}\right) \exp\left(-\frac{r}{L_{\mathbf{B}}}\right), \quad (32)$$

163 where r denotes the spatial distance and $L_{\mathbf{B}}$ is the correlation scale length, fixed as $L_{\mathbf{B}} = 4$
 164 in this application. Being part of Matern kernels, the SOAR function is often used in DA for
 165 prior error correlation modeling [18],[3] thanks to its smoothness and good conditioning. The
 166 simulation of $\mathbf{x}_{b,t} = [u_{b,t}, v_{b,t}]$ via the same discretization of Eq. 30 (except the initial conditions)
 167 is used as background states at time t in the DA modeling. For the knowledge of the exact¹
 168 background error covariance $\mathbf{B}_{E,t}$ at different time, 10^3 background trajectories $\{\mathbf{x}_{b,t}^{\gamma=1 \dots 1000}\}$ are
 169 independently generated via Eq. 31. This exact matrix, hidden for compression approaches,
 170 is seldomly used to evaluate the performance of DA algorithms with reduced observation. To
 171 simulate an industrial context, only 10 trajectories $\{\mathbf{x}_{b,t}^{\gamma=1 \dots 10}\}$ are used in the piecewise esti-
 172 mation of \mathbf{HBH}^T of a flow-independent window, making the ensemble size (10) much smaller
 173 than the problem dimension ($20 \times 20 = 400$).

174 The observations in these twin experiments are generated from the model equivalent based
 175 on the true states (i.e $\mathbf{H}(\mathbf{x}_{\text{true}})$), separately for the fields u and v , respectively denoted as \mathbf{y}_u
 176 and \mathbf{y}_v . For both fields, the observation $\mathbf{y}_t = [\mathbf{y}_{u,t}, \mathbf{y}_{v,t}]$ at time t is the sum of u_t and v_t in a
 177 2×2 cells area with an observation error ϵ_{y_t} ,

$$\mathbf{y}_{u,i,j,t} = u_{\text{true},2i,2j,t} + u_{\text{true},2i+1,2j,t} + u_{\text{true},2i,2j+1,t} + u_{\text{true},2i+1,2j+1,t} + \epsilon_{y_{u,i,j,t}} \quad (33)$$

178 and identical for $\mathbf{y}_{v,i,j,t}$. Thus \mathbf{y} represents also the evolution of the velocity field u and v
 179 with a "coarser" measure as shown in Fig. 1 [g-h].

180 In these experiments, we have set a non-homogeneous observation error covariance where the
 181 error deviation in the center (of radius 4) of the field is 4 times higher, compared to boundary
 182 observations as show in Fig. 3[a]. They are both of the same order of magnitude as $\sigma_{b,0}$,
 183 following also the SOAR function with a smaller scale length $L_{\mathbf{R}} = 1$, compared to background
 184 error correlation. The full error covariance \mathbf{R} of observations \mathbf{y} (after being converted to a 1D
 185 vector by concatenating rows of the original 2D grid model), supposed invariant against time,
 186 is illustrated in Fig. 3 [b]. The \mathbf{R} matrix is supposed to be known in this application, thus
 187 only 10 observation trajectories $\{\mathbf{y}_t^{\gamma=1 \dots 10}\}$ are generated to simulate an ensemble of small
 188 size while evaluating \mathbf{HBH}^T through Eq. 27. In this experiment, we make the choice to
 189 circumvent the difficulty by setting a temporal correlated $\epsilon_{b,t}$, as the background noises are
 190 only added at the beginning of the simulation, and a temporal uncorrelated $\epsilon_{y,t}$. In fact,
 191 temporally correlated background errors are difficult to handle for the Desroziers method since
 192 it treats the innovation quantities as independent samples for covariance estimation. These
 193 assumptions are realistic and widely adopted in DA problems since background simulations
 194 are often taken successively while observations are usually discrete. However, it is beneficial
 195 to have both time uncorrelated $\epsilon_{b,t}$ and $\epsilon_{y,t}$ for Desroziers-type estimation, as long as the error
 196 covariance could still be considered flow-independent [22].

¹Here, by the term "exact", we refer to the covariance truly corresponding to the prior errors present in the background state, no matter the level of optimality of the chosen assimilation scheme.

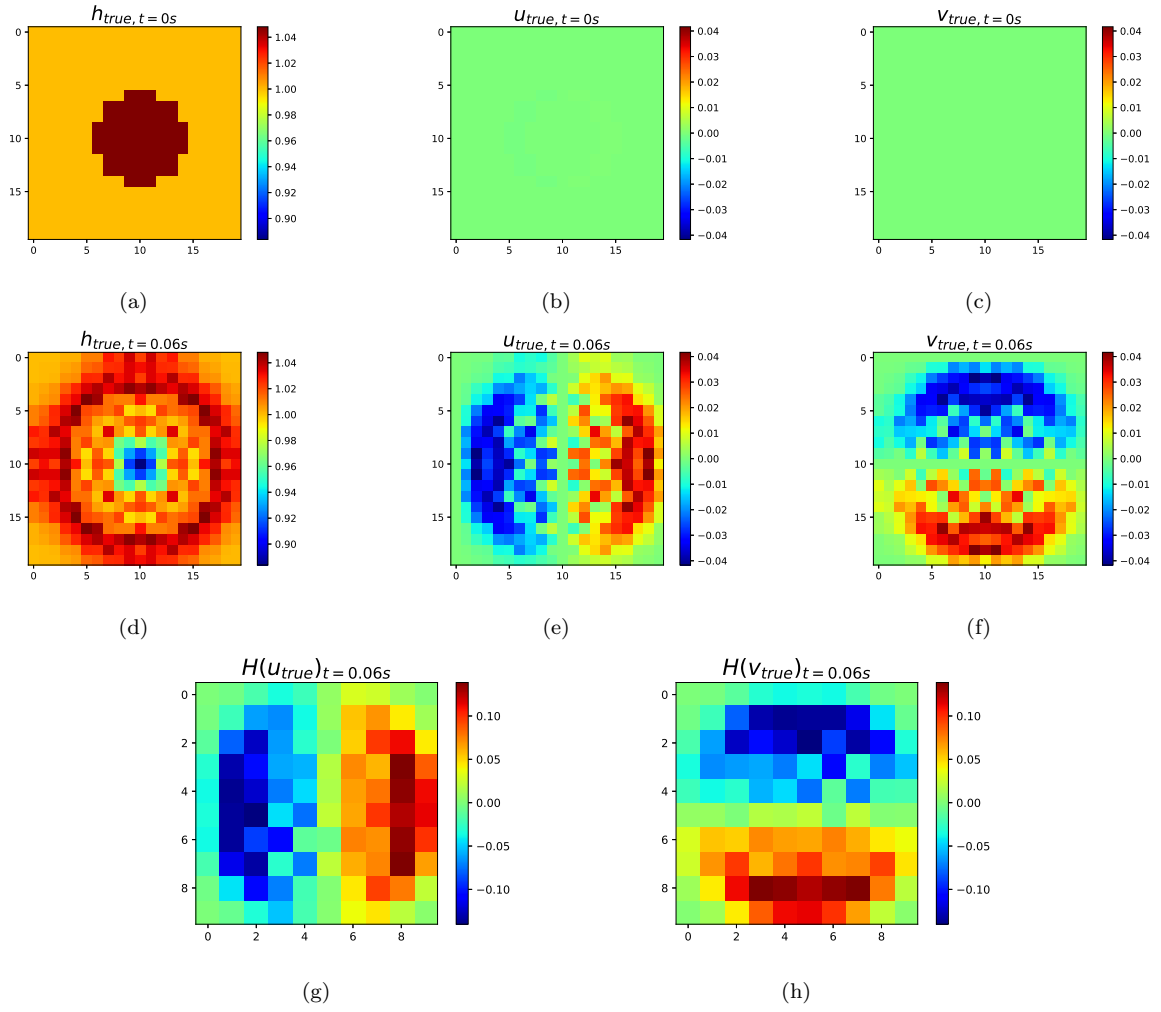


Figure 1: Evolution of the shallow water model of h, u, v (true states) at different time steps (a-f) and the error-free model equivalent $\mathbf{H}(\mathbf{x}_{true})$ for observations (g-h).

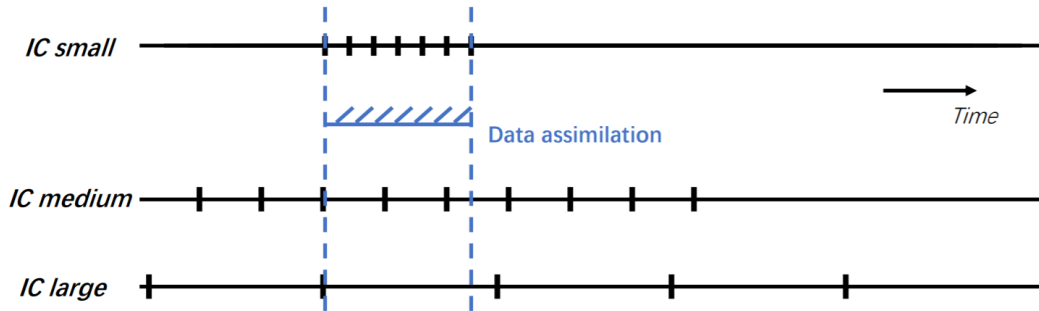


Figure 2: Simple sketch illustrating the three IC sampling strategies. The two vertical blue lines indicate where data assimilation experiments take place (as mentioned in Eq. 34).

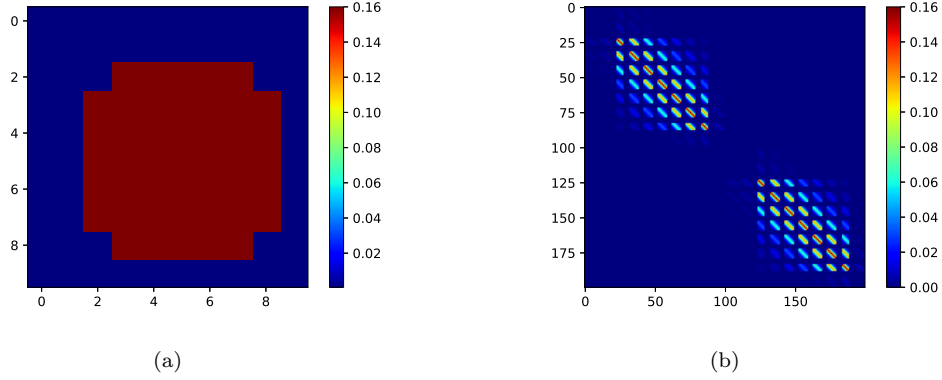


Figure 3: The observation error variance of \mathbf{y}_u and \mathbf{y}_v in the shallow water model[a] and the Balgovind error covariances (\mathbf{R}) after the observation vector (originally in a 2D grid) being converted to a 1D vector [b]

5.2. Numerical results for different compression strategies

We then apply different strategies of observation compression and compare the performance of 3D-Var data assimilation using the reduced observation data. For each assimilation, only the current observation \mathbf{y}_t is used to correct the background state $\mathbf{x}_{b,t}$. Thanks to the 1000 background trajectories $\{\mathbf{x}_{b,t}^{\gamma=1\dots 1000}\}$ simulated, the exact $\mathbf{B}_{E,t}$ matrix can be empirically estimated at different time steps, allowing an accurate estimation of analysis error covariance \mathbf{A}_t via Eq. 8 since the \mathbf{R} matrix is supposed to be known. The matrix trace $\text{Tr}(\mathbf{A})$ then represents the sum of marginal analysis error, equivalent to the square of L^2 norm, i.e $\mathbb{E}(\|\mathbf{x}_a - \mathbf{x}_t\|_2^2)$, often used as an important indicator of DA schemes [18].

Another objective of this experiment is to inspect the impact on the assimilation error given by different sampling densities, which is critical in information-based compression, as stated in the introduction. We display three sampling strategies for \mathbf{HBH}^T estimation with different assumed flow-independent periods $[T_s, T_f]$,

- *IC small*: Dense sampling in a small period, $\Delta_t = 0.001s$ with $T_s = 0.16s$ and $T_f = 0.18s$
- *IC large*: Sparse sampling in a long period, $\Delta_t = 0.1s$ with $T_s = 0s$ and $T_f = 2s$
- *IC medium*: Between *IC small* and *IC large*, with $\Delta_t = 0.01s$ $T_s = 0.1s$ and $T_f = 0.3s$,

as shown in Fig. 2, where Δ_t is the uniform time discretization between two snapshots. For all these three strategies, the \mathbf{HBH}^T is estimated via 20 time steps (i.e $T_f - T_s = 20\Delta_t$), each with 10 background ($\{\mathbf{x}_{b,t}^{\gamma=1\dots 10}\}$) and observation ($\{\mathbf{y}_t^{\gamma=1\dots 10}\}$) states/residuals. To gain a robust comparison, the posterior error variance $\mathcal{E}_{\text{posterior}}$ is averaged using $\text{Tr}(\mathbf{A}_t)$ at four different time, included in all three assumed flow-independent windows,

$$\mathcal{E}_{\text{posterior}} = \frac{\sum \text{Tr}(\mathbf{A}_t)}{4} \quad \text{for } t \in \{0.16, 0.165, 0.170, 0.175\}. \quad (34)$$

We illustrate in Fig. 4[d], the evolution of $\mathcal{E}_{\text{posterior}}$ against the truncation parameter q , varying from 0 to 200. In fact, when $q = 200$, all methods are equivalent since we work with the full observation data. From Fig. 4[d], we observe that all the information-based strategies with different sampling densities are always more optimal compared to the observation-based method for $q \in (0, 200)$. We apply the stopping criteria as described in section 3.3, by calculating the eigenvalues of \mathbf{HBH}^T for the medium sampling strategy. We obtain the optimal truncation parameter $q_{\text{optimal}} = 29$. The distribution of these eigenvalues are shown by the right vertical axes in Fig. 4[d](numerical log scale is represented by the right vertical axis in matching color). With 29 modes, the assimilation correction is achieved from 53.3% to 69.8%, compared to the background model equivalent $\mathcal{H}(\mathbf{x}_b)$ as shown in table 1, which is compatible to the results obtained in [8] when $L_B > L_R$. Among the three sampling strategies, the one of "IC medium"

owns the lowest output error variances, close to the optimal information-based compression where the \mathbf{HBH}^T is computed directly using $\mathbf{B}_{E,t}$. The latter, drawn with blue color in Fig. 4[d], stands for an optimal target for all information-based approaches since we suppose the exact background matrix is out of reach for data compression. As shown in this experiment, the choice of sampling strategy can significantly impact the compression optimality. If the samplings are too close, the residuals might not be uncorrelated, and if the samplings are too sparse, the flow independence of the \mathbf{B} matrix could be threatened. **We remind that the stopping criteria for the truncation parameter q varies for the different sampling strategies as shown in table 1. However, in this experiment the values of the optimal truncation parameters obtained do not qualitatively change the results as shown in Fig. 4[d].**

In Fig. 4[a,c], we display the evolution of the exact background error variances (i.e. $\text{Tr}(\mathbf{B}_{E,t})$) and error correlation (for fixed distances, $r = 1$ and $r = 2$) against time. The estimation of background error correlation in the 2D space, also based on $\mathbf{B}_{E,t}$, is calibrated using the same method shown in [18]. We observe that the error variances increase continuously for both u and v while the spatial error correlation tends to shrink, both being significantly time-variant between $t = 0s$ and $t = 1.4s$. In order to illustrate the non-linear and turbulent nature of error propagation, we show in Fig. 4[b] the error evolution $\|\mathbf{x}_{b,t} - \mathbf{x}_{\text{true},t}\|_2$ of a single background trajectory. Obviously, these facts lead to problem of flow independent assumption for the *IC large* approach (between $0s$ and $2s$), conducting a less optimal compression strategy as shown in Fig. 4[d]. From this twin experiment, we notice the advantage of information-based compression by selecting the most impacting observation components. The optimal sampling strategy may strongly depend on the characteristics (e.g chaosity, stability) of the dynamical system.

Until now, we have shown that, in the idealised case where the observation matrix is known *a priori* and the transformation operator is time-invariant, the information-based approach exhibits advantageous performance compared to the observation-based approach. However, as pointed out by [8], IC approach can be sensitive to prior errors of covariance estimation. In order to investigate the impact of a potential misknowledge of matrix \mathbf{R} , we present here two cases where the difference between the assumed/estimated matrix \mathbf{R}_A and the exact matrix \mathbf{R} is voluntarily large. We explore two cases where the amplitude and the structure are misspecified, respectively:

- (a): \mathbf{R}_A has the same correlation structure as \mathbf{R} with an homogeneous marginal error variance (i.e. $\mathbf{R}_{A,i,i} = 0.04$ which is different to \mathbf{R} (cf. Fig. 3(a)).)
- (b): The correlation scale L_{R_A} is set to be 5 while $L_R = 1$ as explained in section 5.1 with same marginal error variances.

In both cases, the observation compression is implemented using \mathbf{R}_A while the observation matrix in the reduced space is set to be $\mathbf{R}_A^{-1/2} \mathbf{R} \mathbf{R}_A^{-1/2}$ instead of the identity matrix in Eq. 16 and Eq. 22. The performance of these compression methods is illustrated in Fig. 5, respectively for case (a) and (b). The optimal IC solutions (same as the blue lines in Fig. 4(d)) are drawn in dashed blue lines for comparison purposes. Both OC and IC approaches exhibit less optimal performance compared to Fig. 4. Furthermore, as shown in Fig. 5(a), the IC method can be more sensitive to the mis-specification of the \mathbf{R} matrix amplitude, leading in this case to larger output error variances while IC behaves better than OC for misspecified \mathbf{R} matrix correlation length.

	OC	IC large	IC medium	IC small	IC optimal
q_{optimal}	22	48	29	25	78
Correction for $q = 29$	53.3%	61.5%	65.7%	62.7%	69.8%

Table 1: The ratio of background minus analysis innovation ($\|\mathcal{H}(\mathbf{x}_b) - \mathcal{H}(\mathbf{x}_a)\|_2$) using compressed observation, relative to the one obtained with full observation

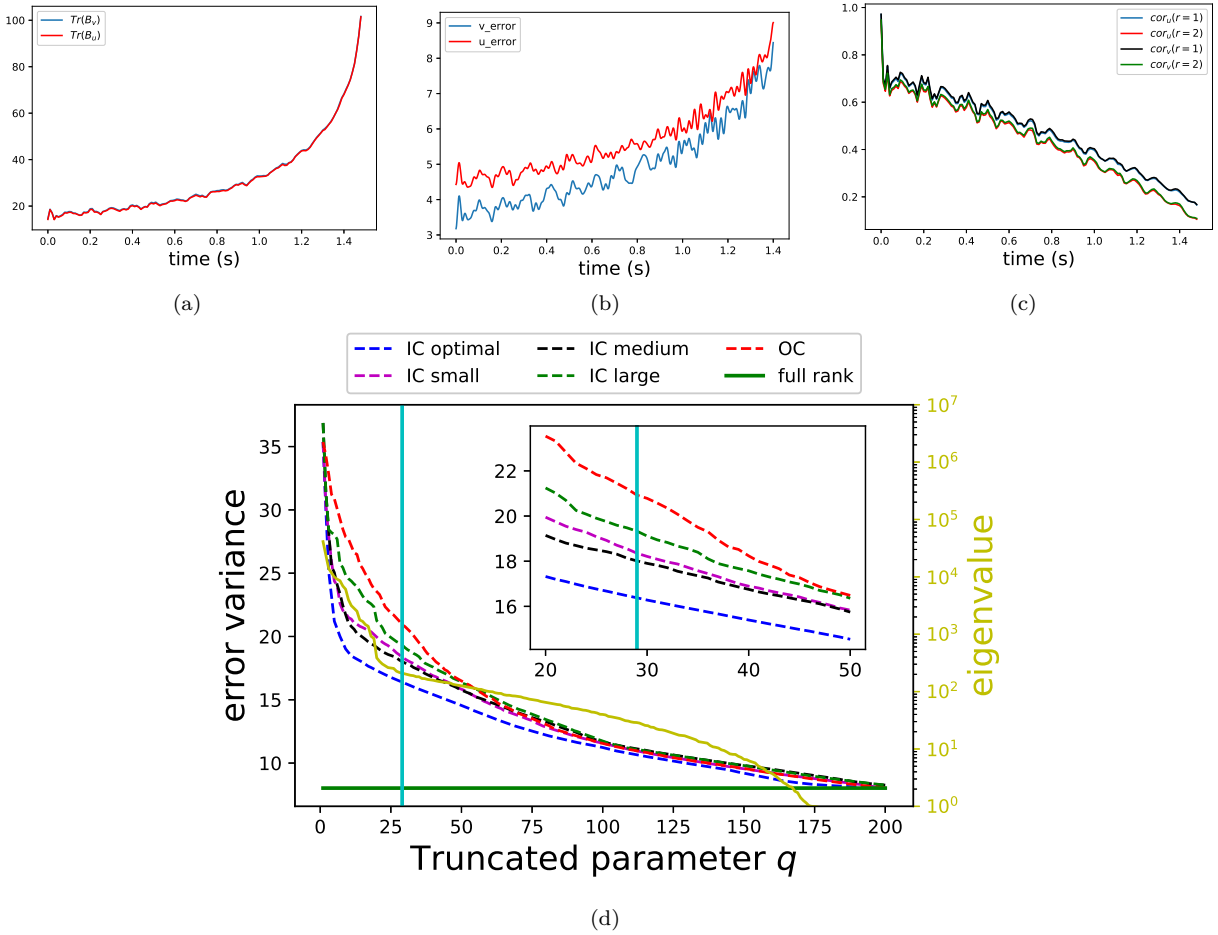


Figure 4: [a]: evolution of the exact background variance of u ($\text{Tr}(\mathbf{B}_u)$) and v ($\text{Tr}(\mathbf{B}_v)$) against time; [b]: evolution of $\|\mathbf{x}_{b,t} - \mathbf{x}_{\text{true},t}\|_2$ of a single background trajectory; [c]: evolution of average error correlation, of fixed distances ($r = 1$ and $r = 2$) in the 2D space; [d]: analysis error variance $\mathcal{E}_{\text{posterior}}$ (left y-axis) and eigenvalues of the estimated background error covariance in observation space ($\mathbf{H}\mathbf{B}\mathbf{H}^T$) (right y-axis) of the medium sampling strategy as a function of the truncation parameter for $t \in [0.16s, 0.18s]$. The vertical line represents the stopping criteria of $\sigma_q = \sqrt{\sigma_1}$;

6. Application to an operating hydrological model

6.1. DA modelling for flow reanalysis/prediction

The compression strategies introduced in previous sections are applied to a hydrological application using a precipitation-flow simulator MORDOR-TS developed by Électricité de France (EDF, the French electric utility company). This software is widely applied in operating hydraulic/hydrological problems, e.g. [39], [24], [40]. Based on information on spatially distributed physical parameters, such as precipitation or temperature, it provides a simulation of river flow relying on conceptual watersheds modeling. For more details about MORDOR-TS, interested readers are referred to [24] and [23]. MORDOR-TS is used as a non-linear state-observation transformation operator in data assimilation. We concentrate on a study area in the south of France, around the Tarn river where 9 streamflow gauges positioned at different mesh outlets are available. The Tarn river, being known for its extreme variability of water-level values and high sensitivity to precipitations [23], is an ideal benchmark for comparing different DA strategies. Located downstream, the Tarn river outlet at Millau (hereby denoted as TM) is of particular interest in the hydrological study. As an example, we show in Fig. 6 the simulated and daily observed Tarn river discharges at Millau, for 3 months in 1990 with the averaged precipitation over 28 spatially distributed regions (see [23]). Significant impacts of precipitation on the river flow of TM is observed with a delay of 2 to 5 days. The objective of this DA modeling is to improve the river flow prediction and reanalysis (history matching) by performing corrections on the daily precipitation in the 28 regions. Other physical quantities (e.g temperature) are

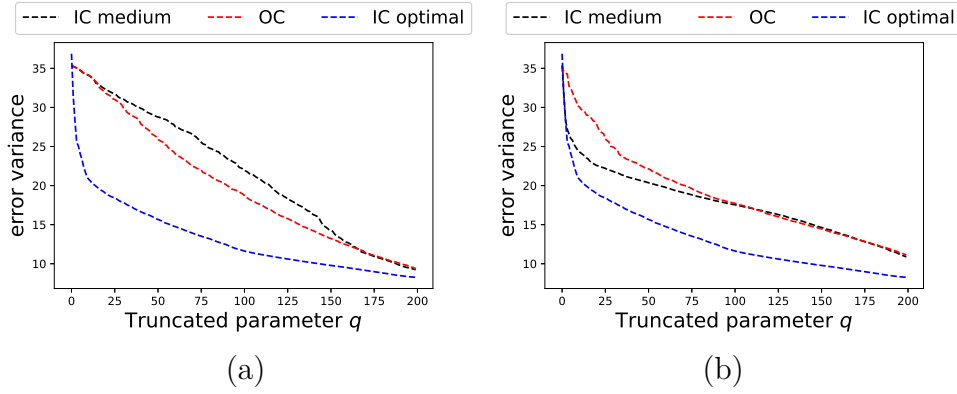


Figure 5: Evolution of error variance when the observation matrix is mis-specified.

288 considered as invariant parameters in this study. The variational assimilation is performed
 289 using the ADAO [41] package of SALOME platform, also developed by EDF.

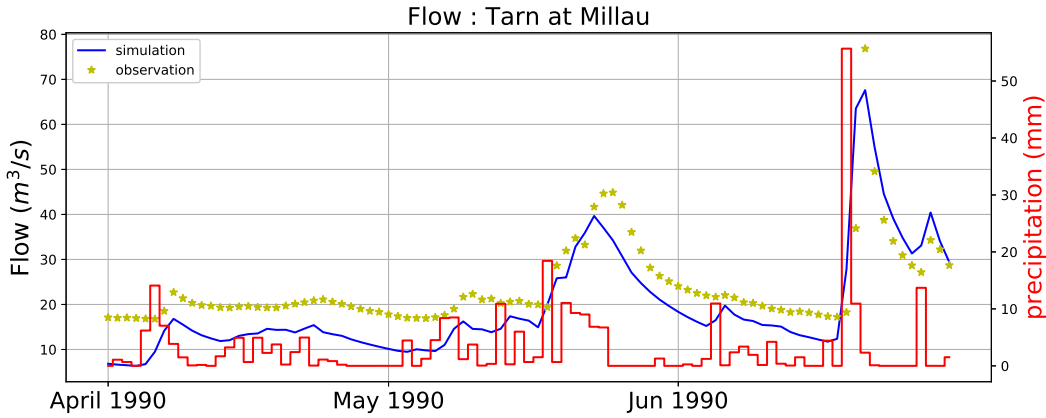


Figure 6: Example of simulation predicted by MORDOR-TS using daily precipitation, and observed Tarn discharges at Millau for three months in 1990. Simultaneous observed precipitations are in red bars (with the scale on the right vertical axis).

290 As mentioned in [23], performing DA correction on all precipitation inputs (i.e 28 regions)
 291 can probably introduce an over-parameterization and thus induces an overfitting, with a high
 292 risk to deteriorate flow forecasts. Therefore, we make the choice to proceed with uniform
 293 additional increments ξ_t^p for all 28 regions, depending only on time t . **Incremental variables**
 294 $\xi^{r,j}$ ($j = 1..8$) on the eight parameters which determine the initial (at $t = 0$) reservoir level
 295 is also added **in the state space** to adjust the river flow at the beginning of each assimilation
 296 window. These windows are fixed of 30 days, leading to an observation vector of dimension
 297 270 with 9 gauges. Temporal correlation is considered for both background and observation
 298 errors. The DA modelling is summarized in Table 2 and a more detailed description can be
 299 found in [23] and [40]. The main objective of this application stands for improving short-range
 300 flow forecasting by correcting historical precipitation. Since the impact of the precipitation on
 301 the river flow is only significant within 3 to 4 days (see [23] for details), we fix the prediction
 302 window to 3 days in this study.

DA modelling	state: \mathbf{x}	$dim(\mathbf{x})$	Observations: \mathbf{y}	$dim(\mathbf{y})$	invariant parameters
Incremental 3DVar	ξ_t^p $\xi^{r,j}$	38	river flow $Q_{q,t}$	270	temperature, etc

Table 2: Details of DA modelling where $t = 0..29$ is the time (days) relative to the beginning of the assimilation window; $q = 1..9$ represents the 9 gauges where $\xi_t^p, \xi^{r,j}$ represent respectively the increments of daily precipitation and initial reservoir level.

303 6.2. Observation compression

304 Despite that MORDOR-TS is computationally efficient (it may take only a few CPU seconds
 305 to simulate a spatially distributed flow simulation of several years), the application of variational
 306 assimilation algorithms could be expensive, due to the non-linearity of the transformation
 307 operator. As shown in Table. 2, in this DA modeling, the dimension of the observation vector
 308 is much larger compared to the state dimension, promoting the utilization of observation data
 309 compression. We then implement DA algorithms in the hydrological model with compressed
 310 data using either OC or IC approaches. **To make the compression strategies more general**, in
 311 both cases the principal components are constructed using the daily observed flow data from
 312 1990 to 2000 in the 9 gauges. The objective of this study is to make an efficient use of the
 313 observation vector \mathbf{y} with an optimal number of modes selected, which we expect to be much
 314 smaller than the full observation dimension ($\dim(\mathbf{y}) = 270$).

A major hurdle of this application is that the *a priori* knowledge of both \mathbf{B} and \mathbf{R} is very limited. As a remedy, we start as described in [23], by considering the background covariance matrix \mathbf{B} of Balgovind-type since we wish to model the existence of temporal correlation in the precipitation data. Moreover, the initial \mathbf{R} matrix is set to be diagonal. The DI01 algorithm [20] is then applied several times to come up with a reasonable approximation of the ratio between $\text{Tr}(\mathbf{B})$ and $\text{Tr}(\mathbf{R})$ at the first stage. In a second stage, we then perform the estimation of \mathbf{HBH}^T and \mathbf{R} , relying on Desroziers formulation (respectively Eq. 26 and 28) using 3400 assimilation windows of 30 days from 1990 to 2000. By then, post-processing is required to ensure the symmetric positive definiteness (SPD) of the \mathbf{R} matrix. More precisely,

$$\mathbf{R} \leftarrow \frac{1}{2}(1 - \mu)(\mathbf{R} + \mathbf{R}^T) + \mu\mathbf{C}, \quad (35)$$

315 where $\mu = 0.1$ and $\mathbf{C} = \text{Tr}(\mathbf{R}) \times \mathbf{I}$. The Desroziers method is iterated twice, using the same
 316 data set, to ensure the stability of the estimated matrices. The algorithm outputs produced
 317 after the first and the second iterations are very similar as shown in [23]. We emphasize that
 318 the estimated \mathbf{R} matrix is not only used for the observation compression but also in the DA
 319 algorithm in the full observation space. The \mathbf{HBH}^T matrix is obtained through Eq. 27, once
 320 the \mathbf{R} matrix is specified. As a remark, even if the system considered here is not very large, the
 321 computational burden associated with the data assimilation of this nonlinear system (for which
 322 prior information is degraded) remains important because of a multi-stage tuning approach
 323 which combined several offline and online covariance tuning algorithms can be implemented to
 324 improve the reanalysis and the forecasting accuracy of this hydrological application. However,
 325 these methods are computationally expensive, especially when iterations are needed (e.g. [18]).
 326 With advanced data compression methods, the computational burden can be released, allowing
 327 more precise covariance tuning to improve the DA performance.

328 6.3. DA with compressed data

329 6.3.1. Averaged performance

Extracting the principal components $\tilde{\mathbf{L}}$ and $\hat{\mathbf{L}}$, respectively based on estimated \mathbf{HBH}^T and \mathbf{R} , we then apply the compression methodology described in sect.3. The objective is to compare the assimilation output $\mathbf{x}_{a,\text{compression}}$ and $\mathbf{x}_{a,\text{full}}$, obtained using either the compressed observation $\hat{\mathbf{y}}_q$, $\tilde{\mathbf{y}}_q$ or the full observation vector \mathbf{y} . More precisely, we are interested in the observation minus analysis (O-A) innovation quantity for both flow reanalysis and forecast. Varying the truncated parameter q , DA processes are performed respectively with $(\mathbf{x}_b, \tilde{\mathbf{y}}_q, \mathbf{B}, \tilde{\mathbf{R}}_q, \tilde{\mathcal{H}}_q)$ and $(\mathbf{x}_b, \hat{\mathbf{y}}_q, \mathbf{B}, \hat{\mathbf{R}}_q, \hat{\mathcal{H}}_q)$ for 12 assimilation windows in 1993, each of 30 days starting at the first day of every month. We draw the averaged compressed/full O-A innovation ratio ∇ , defined as

$$\nabla = \frac{\|\mathbf{y} - \mathcal{H}(\mathbf{x}_{a,\text{compression}})\|_2}{\|\mathbf{y} - \mathcal{H}(\mathbf{x}_{a,\text{full}})\|_2}, \quad (36)$$

330 in Fig. 7 for both reanalysis[a] and prediction[b] at TM. More particularly, $\nabla = 100\%$ means
 331 the reanalysis/prediction accuracy of the current solution is equivalent to the one obtained with
 332 the full observation vector.

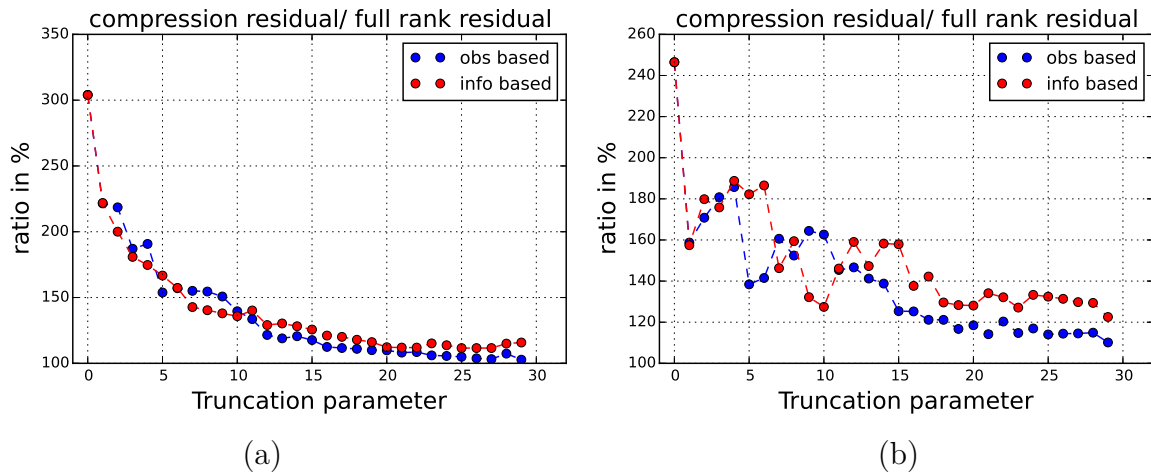


Figure 7: Evolution of ∇ , averaged using 12 assimilation windows, against the number of truncation parameter q for reanalysis[a] and prediction[b] at TM.

333 We observe from Fig. 7 that the performance of these two approaches is similar to the
 334 reanalysis while the observation-based method is slightly more optimal on average for flow fore-
 335 casting. The evolution of the reconstruction error (Fig. 7[a]) is much smoother, compared to the
 336 prediction error ((Fig. 7[b])), both against the truncation parameter. In fact, the reconstruction
 337 error is estimated using assimilation windows of 30 days while prediction windows are solely of
 338 3 days. Therefore, the estimation of the prediction ratio has significantly more sampling noise.
 339 Furthermore, since both \mathbf{B} and \mathbf{R} are not well specified *a priori*, extra noise can be introduced
 340 while estimating the information entropy. For both methods, the assimilation results obtained
 341 using 15 to 20 modes (around 5% to 7.5% of total observation dimension) are close to the full
 342 rank solution in terms of both reanalysis and prediction. Without deteriorating the assimilation
 343 result, these compression strategies make the DA algorithm certainly more efficient, allowing
 344 more optimization iterations if needed.

345 6.3.2. Performance in each DA window

346 We draw the reconstructed river flow (i.e $\mathcal{H}(\mathbf{x}_a)$) at TM of each of those 12 assimilation
 347 windows, for both corrections with compressed and full observation data in Fig. 8 where the
 348 yellow stars represent the daily observations. Based on the method described in Eq. 25, the
 349 optimal truncation parameter reads $q_{\text{optimal}} = 22$. Here, we display the results when $q = 10$
 350 in order to voluntarily emphasize the difference between the two approaches as shown in Fig 7.
 351 A vertical line in each graph separates the reanalysis (left) and the prediction (right). We
 352 notice that the reconstructed curves issued from OC (blue) and IC (red) are similar in most
 353 cases, both being adequately close to the full rank assimilation (green), compared to the original
 354 simulation. Some exceptions can be found, for example, in the assimilation window of December
 355 1993 where the prediction is covered by a flood period. It seems that the information-based
 356 approach provides a better performance, especially for flow forecasting at that moment. In
 357 general, as demonstrated in [23], meteorological factors can impact the assimilation precision
 358 significantly. DA algorithms often perform better during drought periods (see Jun, July, August
 359 in Fig. 8) where the prior observation minus background (O-B) innovations are more consistent
 360 (i.e always being under-estimated or over-estimated). Contrarily, in flood periods where O-B
 361 innovations are usually more turbulent, more careful attention might be taken when performing
 362 compression methods.

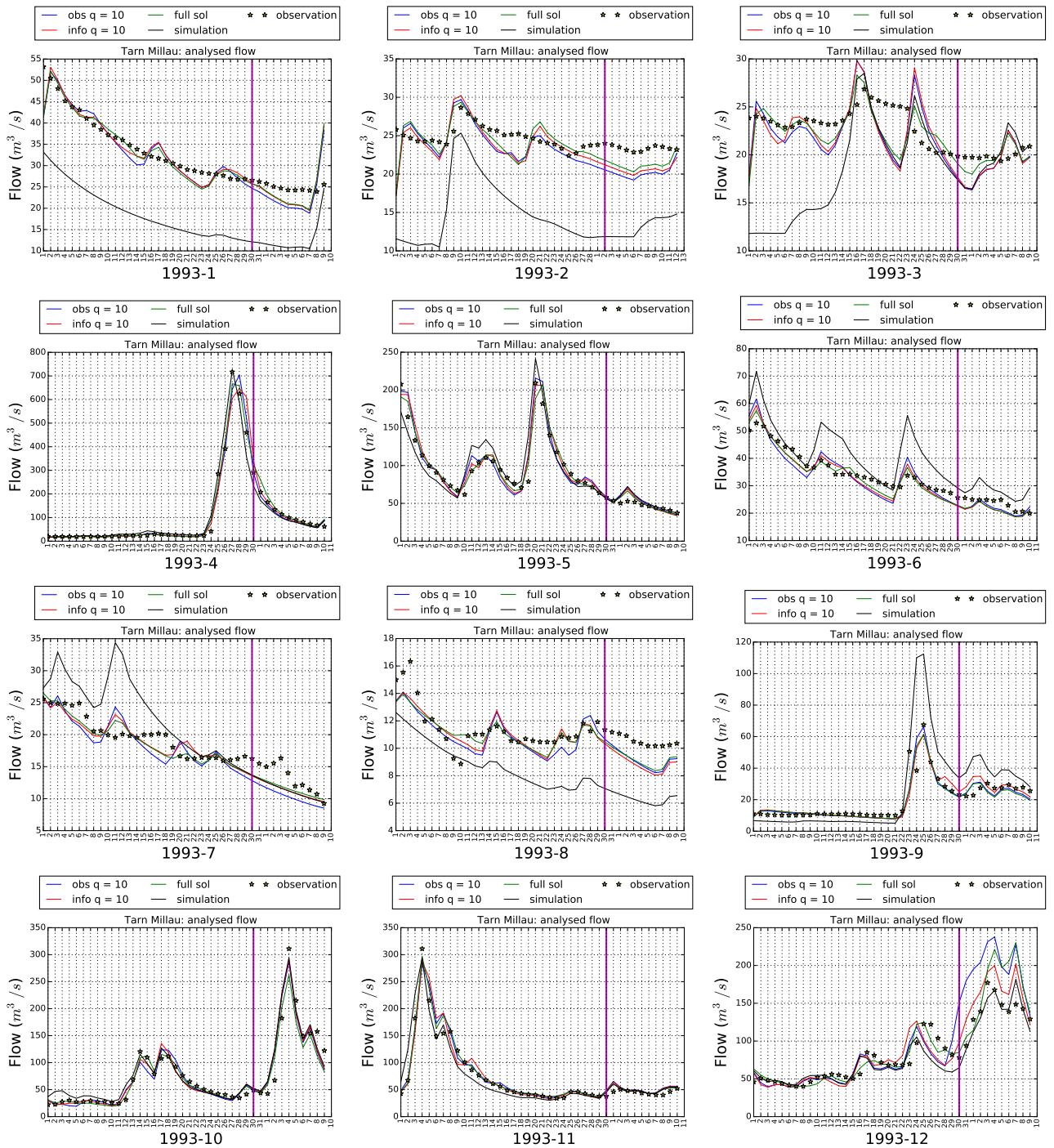


Figure 8: The reconstructed and predicted river flow at TM for OC (obs) ($q = 10$), IC (info) ($q = 10$) and full rank DA solutions of different months in 1993. The left side of the vertical line represents the flow reanalysis while the right side represents the prediction

363 7. Discussion

364 Sequential data assimilation algorithms can be computationally challenging, especially for
365 large scale systems such as NWP, remote sensing, or geophysical problems. Data compression
366 techniques commonly used in DA problems have recently received increasing interest in reducing
367 the computational burden. Much effort has been devoted to improving the algorithm efficiency
368 without diminishing the accuracy of assimilation reconstruction and forecasting. Classical com-
369 pression approaches consist of either extracting the principal vectors of observation dynamics or
370 identifying the directions that contribute the most to the prior-posterior information gap. For
371 both methods, the lack of precise knowledge on prior error covariances stands for an essential
372 obstacle, as mentioned in several previous studies. Furthermore, the limited number of back-
373 ground/observation trajectories often entails a poor empirical estimation. In this paper, we have
374 introduced a concept of observation compression benefiting from existing piecewise covariance
375 estimation, establishing a natural connection between the posterior error covariance diagnosis
376 and data compression techniques. More precisely, we assume that the error covariances (both \mathbf{B}
377 and \mathbf{R}) are flow-independent over some specific time periods, which allows an estimation based
378 on time-variant residuals. Therefore, a much smaller number of background/observation trajec-
379 tories are required for non-parametric covariance estimation. Different estimation formulations
380 are possible depending on the prior knowledge of the \mathbf{R} matrix. The choice of flow-independent
381 windows, as well as the residual sampling densities, is essential in these approaches, especially
382 for the \mathbf{HBH}^T estimation. When the samplings are either too dense or too sparse, the as-
383 sumptions of covariance estimation approaches might be unsatisfied, leading to a less optimal
384 observation compression. These aspects are numerically analyzed in the twin experiments of a
385 2D shallow water model with non-linear dynamics with the perfect knowledge of the \mathbf{R} . Nu-
386 merical results show a significant advantage of the information-based compression in terms of
387 assimilation accuracy, compared to the observation-based one. As for the industrial hydrologi-
388 cal model, posterior covariance estimation which requires the knowledge of the analyzed states
389 \mathbf{x}_a , is needed since the \mathbf{R} matrix is not known *a priori*. In this application, both the OC and
390 IC compression methods rely on the flow-independent estimation of \mathbf{R} , showing competitive
391 performance regarding the flow reanalysis and the forecasting accuracy. A meteorological effect
392 is also briefly discussed in this hydrological application, which indicates that different numbers
393 of modes should be chosen in different periods of the year regarding the hydrological proper-
394 ties. Future work can be considered to improve the algorithm efficiency and flexibility under
395 industrial conditions, for example, by using parametric covariance tuning methods or spatial
396 localization techniques. Another important limitation of the current approach stands for the
397 time invariance of the observations error covariance on some time-scale, limiting for instance
398 the use of moving observation sensors. Indeed, if observation positions change, both the obser-
399 vation matrix \mathbf{R} and the transformation operator \mathcal{H} can not be considered flow-independent,
400 leading to difficulties when applying Desroziers-type methods. Future work can be considered
401 to use interpolation approaches to construct a global observation set which includes all time-
402 variant observation positions. Another perspective of this study could be to further examine
403 the optimal choice of the sampling density while estimating the error covariances, for example,
404 with the help of uncertainty quantification methods for dynamical systems.

405 Acknowledgement

406 The authors would like to thank Dr. Bertrand Iooss and Dr. Angélique Ponçot for fruitful
407 discussions about the compression methodology and the hydrological application. This work
408 was supported by EDF R&D. This research was partially funded by the Leverhulme Centre for
409 Wildfires, Environment and Society through the Leverhulme Trust, grant number RC-2018-023.
410 The authors are grateful to an anonymous reviewer for the useful remarks on the manuscript.

411 Bibliography

- 412 [1] F. Rabier, Overview of global data assimilation developments in numerical weather-
413 prediction centres, *Quarterly Journal of the Royal Meteorological Society* 131 (2005)
414 3215–3233.
- 415 [2] M. C. Rochoux, S. Ricci, D. Lucor, B. Cuenot, A. Trouvé, Towards predictive data-driven
416 simulations of wildfire spread—part i: Reduced-cost ensemble kalman filter based on a
417 polynomial chaos surrogate model for parameter estimation, *Natural Hazards and Earth*
418 *System Sciences* 14 (2014) 2951–2973.
- 419 [3] H. Gong, Y. Yu, Q. Li, C. Quan, An inverse-distance-based fitting term for 3D-Var data
420 assimilation in nuclear core simulation, *Annals of Nuclear Energy* 141 (2020) 107346.
- 421 [4] S. Cheng, R. Arcucci, C. C. Pain, Y.-K. Guo, Optimal vaccination strategies for
422 covid-19 based on dynamical social networks with real-time updating, *arXiv preprint*
423 *arXiv:2103.00485* (2021).
- 424 [5] D. Lucor, O. P. Le Maître, Cardiovascular modeling with adapted parametric inference,
425 *ESAIM: ProcS* 62 (2018) 91–107.
- 426 [6] P. Nadler, R. Arcucci, Y. Guo, Data assimilation for parameter estimation in economic
427 modelling, in: *15th International Conference on Signal-Image Technology & Internet-Based*
428 *Systems (SITIS) 2019*, 2019, pp. 649–656.
- 429 [7] A. D. Collard, A. P. McNally, F. I. Hilton, S. B. Healy, N. C. Atkinson, The use of principal
430 component analysis for the assimilation of high-resolution infrared sounder observations
431 for numerical weather prediction, *Quarterly Journal of the Royal Meteorological Society*
432 136 (2010) 2038–2050.
- 433 [8] A. Fowler, Data compression in the presence of observational error correlations, *Tellus A:*
434 *Dynamic Meteorology and Oceanography* 71 (2019) 1634937.
- 435 [9] A. Carrassi, M. Bocquet, L. Bertino, G. Evensen, Data assimilation in the geosciences: An
436 overview of methods, issues, and perspectives, *Wiley Interdisciplinary Reviews: Climate*
437 *Change* 9 (2018) e535.
- 438 [10] S. Cheng, J.-P. Argaud, B. Iooss, A. Ponçot, D. Lucor, A graph clustering approach to
439 localization for adaptive covariance tuning in data assimilation based on state-observation
440 mapping, accepted for publication in *Mathematical Geosciences*, 2021. [arXiv:2001.11860](https://arxiv.org/abs/2001.11860).
- 441 [11] A. Cioaca, A. Sandu, Low-rank approximations for computing observation impact in 4D-
442 Var data assimilation, *Computers & Mathematics with Applications* 67 (2014) 2112 –
443 2126.
- 444 [12] E. D. Nino-Ruiz, A. Sandu, X. Deng, A parallel implementation of the ensemble kalman
445 filter based on modified cholesky decomposition, *Journal of Computational Science* 36
446 (2019) 100654.
- 447 [13] I. Hoteit, D.-T. Pham, J. Blum, A simplified reduced order Kalman filtering and appli-
448 cation to altimetric data assimilation in tropical pacific, *Journal of Marine Systems* 36
449 (2002).
- 450 [14] J.-P. Argaud, B. Bouriquet, F. Caso, H. Gong, Y. Maday, O. Mula, Sensor placement
451 in nuclear reactors based on the generalized empirical interpolation method, *Journal of*
452 *Computational Physics* 363 (2018) 354 – 370.

- 453 [15] J. A. Waller, S. L. Dance, N. K. Nichols, On diagnosing observation-error statistics with
454 local ensemble data assimilation, *Quarterly Journal of the Royal Meteorological Society*
455 143 (2017) 2677–2686.
- 456 [16] M. Matricardi, A. P. McNally, The direct assimilation of principal components of IASI
457 spectra in the ECMWF 4D-Var, *Quarterly Journal of the Royal Meteorological Society*
458 140 (2014) 573–582.
- 459 [17] M. Fisher, Background error covariance modelling, in: Seminar on Recent developments
460 in data assimilation for atmosphere and ocean (Shinfield Park, Reading, 8-12 September),
461 ECMWF, 2003.
- 462 [18] S. Cheng, J.-P. Argaud, B. Iooss, D. Lucor, A. Ponçot, Background error covariance
463 iterative updating with invariant observation measures for data assimilation, *Stochastic*
464 *Environmental Research and Risk Assessment* 33 (2019) 2033–2051.
- 465 [19] D. F. Parrish, J. C. Derber, The National Meteorological Center’s spectral statistical-
466 interpolation analysis system, *Monthly Weather Review* 120 (1992) 1747–1763.
- 467 [20] G. Desroziers, S. Ivanov, Diagnosis and adaptive tuning of observation-error parameters
468 in a variational assimilation, *Quarterly Journal of the Royal Meteorological Society* 127
469 (2001) 1433 – 1452.
- 470 [21] G. Desroziers, L. Berre, B. Chapnik, P. Poli, Diagnosis of observation, background and
471 analysis-error statistics in observation space, *Quarterly Journal of the Royal Meteorological*
472 *Society* 131 (2005) 3385 – 3396.
- 473 [22] K. Bathmann, Justification for estimating observation-error covariances with the
474 Desroziers diagnostic, *Quarterly Journal of the Royal Meteorological Society* 144 (2018)
475 1965–1974.
- 476 [23] S. Cheng, J.-P. Argaud, B. Iooss, D. Lucor, A. Ponçot, Error covariance tuning in vari-
477 ational data assimilation: application to an operating hydrological model, *Stochastic*
478 *Environmental Research and Risk Assessment* 35 (2021) 1019–1038.
- 479 [24] L. Rouhier, M. Le Lay, F. Garavaglia, N. Moine, F. Hendrickx, C. Monteil, P. Ribstein, Im-
480 pact of mesoscale spatial variability of climatic inputs and parameters on the hydrological
481 response, *Journal of Hydrology* 553 (2017) 13 – 25.
- 482 [25] W. Fulton, Eigenvalues, invariant factors, highest weights, and schubert calculus, *Bulletin*
483 *of The American Mathematical Society* 37 (2000) 209–250.
- 484 [26] F. Uboldi, M. Kamachi, Time-space weak-constraint data assimilation for nonlinear mod-
485 els, *Tellus A* 52 (2000) 412–421.
- 486 [27] J. Brajard, A. Carrassi, M. Bocquet, L. Bertino, Combining data assimilation and machine
487 learning to emulate a dynamical model from sparse and noisy observations: A case study
488 with the Lorenz 96 model, *Journal of Computational Science* 44 (2020) 101171.
- 489 [28] P. Antonelli, H. E. Revercomb, L. A. Sromovsky, W. L. Smith, R. O. Knuteson, D. C.
490 Tobin, R. K. Garcia, H. B. Howell, H.-L. Huang, F. A. Best, A principal component noise
491 filter for high spectral resolution infrared measurements, *Journal of Geophysical Research:*
492 *Atmospheres* 109 (2004).
- 493 [29] D. Tobin, P. Antonelli, H. Revercomb, S. Dutcher, D. Turner, J. Taylor, R. Knuteson,
494 K. Vinson, Hyperspectral data noise characterization using principle component analysis:
495 Application to the atmospheric infrared sounder, *Journal of Applied Remote Sensing* 1
496 (2006) 013515.

- 497 [30] C. Cardinali, S. Pezzulli, E. Andersson, Influence-matrix diagnostic of a data assimilation
498 system, *Quarterly Journal of the Royal Meteorological Society* 130 (2004) 2767–2786.
- 499 [31] S. Migliorini, Information-based data selection for ensemble data assimilation, *Quarterly*
500 *Journal of the Royal Meteorological Society* 139 (2013) 2033–2054.
- 501 [32] R. Cangelosi, A. Goriely, Component retention in principal component analysis with
502 application to cdna microarray data, *Biology Direct* 2 (2007) 2.
- 503 [33] R. Arcucci, L. Mottet, C. Pain, Y.-K. Guo, Optimal reduced space for variational data
504 assimilation, *Journal of Computational Physics* 379 (2018) 51–69.
- 505 [34] C. Lanczos, An iteration method for the solution of the eigenvalue problem of linear
506 differential and integral operators, *Journal of research of the National Bureau of Standards*
507 45 (1950) 255–282.
- 508 [35] P. Tandeo, P. Ailliot, M. Bocquet, A. Carrassi, T. Miyoshi, M. Pulido, Y. Zhen, A review
509 of innovation-based methods to jointly estimate model and observation error covariance
510 matrices in ensemble data assimilation, *Monthly Weather Review* (2020) 1–68.
- 511 [36] A. Hollingsworth, P. Lönnberg, The verification of objective analyses: Diagnostics of
512 analysis system performance, *Meteorology and Atmospheric Physics* 40 (1989) 3–27.
- 513 [37] J. A. Waller, S. L. Dance, N. K. Nichols, Theoretical insight into diagnosing observa-
514 tion error correlations using observation-minus-background and observation-minus-analysis
515 statistics, *Quarterly Journal of the Royal Meteorological Society* 142 (2016) 418–431.
- 516 [38] B. Chapnik, G. Desroziers, F. Rabier, O. Talagrand, Property and first application of an
517 error-statistics tuning method in variational assimilation, *Quarterly Journal of the Royal*
518 *Meteorological Society* 130 (2004) 2253 – 2275.
- 519 [39] R. Garçon, Préviation opérationnelle des apports de la Durance à Serre-Ponçon à l’aide du
520 modèle MORDOR. Bilan de l’année 1994-1995, *La Houille Blanche* (1996) 71–76.
- 521 [40] S. Cheng, Error covariance specification and localization in data assimilation with indus-
522 trial application, Ph.D. thesis, Paris-Saclay University, France, 2020.
- 523 [41] J.-P. Argaud, User documentation, in the SALOME 9.3 platform, of the ADAO module for
524 ”Data Assimilation and Optimization”, Technical report 6125-1106-2019-01935-EN, EDF /
525 R&D, 2019.