



# Real-time estimation and prediction of unsteady flows using reduced-order models coupled with few measurements

Valentin Resseguier, Matheus Ladvig, Dominique Heitz

## ► To cite this version:

Valentin Resseguier, Matheus Ladvig, Dominique Heitz. Real-time estimation and prediction of unsteady flows using reduced-order models coupled with few measurements. *Journal of Computational Physics*, 2022, 471, pp.111631. 10.1016/j.jcp.2022.111631 . hal-03445455v3

**HAL Id: hal-03445455**

**<https://hal.science/hal-03445455v3>**

Submitted on 2 Mar 2022

**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.

# Real-time estimation and prediction of unsteady flows using reduced-order models coupled with few measurements

Valentin Resseguier\*, Matheus Ladvig

*Lab, SCALIAN DS, Espace Nobel, 2 Allée de Becquerel, Rennes, France*

Dominique Heitz

*INRAE, OPAALE, Rennes, France*

---

## Abstract

The estimation and prediction of unsteady flows in real time offers significant advantages for the monitoring and active control of complex hydrodynamic and aerodynamic systems, such as wind turbine blades, hydrofoils and aircraft wings. A new data assimilation algorithm is proposed for the estimation and prediction of unsteady flows, coupling in real time onboard measurements and fluid dynamics simulations at minimal computational expense. The procedure combines a Proper Orthogonal Decomposition Galerkin method, a model under location uncertainty stochastic closure, and a particle filtering scheme. The algorithm is validated using case studies of two- and three-dimensional wake flows at low and moderate Reynolds numbers respectively. Following an initial learning window to train the algorithm, and using only a single measurement point, our method is shown to perform well against conventional reduced data assimilation algorithms for up to 14 vortex shedding cycles.

*Keywords:* Fluid dynamics, reduced order model, uncertainty quantification, stochastic closure, particle filtering

---

## 1. INTRODUCTION

Active control of complex aeroelastic and aerodynamic systems has the potential to yield significant advantages, from the alleviation of load on turbine blades and optimization of electricity production from wind farms [1, 2], to active flutter suppression for aircraft [3]. Such active controls can require state observers but estimating – or even predicting – an unsteady turbulent flow state from sparse measurements in real-time can be challenging. Through statistical estimation techniques, sensor observations can be assimilated to flow dynamical models’ predictions, but some difficulties must first be overcome for this data assimilation method to be a viable strategy.

Firstly, the simulation model is subject to a number of conditions. Notably, the simulation must resolve sufficient spatiotemporal scales for the data assimilation outputs to be stable, specifically in real-time. Purely data-driven fluid

---

\*Corresponding author

Email address: [valentin.resseguier@scaliam.com](mailto:valentin.resseguier@scaliam.com) (Valentin Resseguier)



dynamics models employing, for instance, machine learning techniques have the potential to achieve a significant reduction in computation time. But, they generally require extensive data to assimilate, or else it may not be sufficiently robust for accurate and stable predictions of turbulent flows. As detailed later in this paragraph, intermittency, sensitivity to perturbations and closure issues are probably responsible for these limitations. Time-wise fluid flow estimations that do not consider underlying time dynamics – such as supervised linear and non-linear interpolations [e.g. 4] – are likely to suffer from similar limitations. Conversely, pure physics-based models, such as Large Eddy Simulation (LES) and Reynolds-averaged Navier-Stokes (RANS) remain too computationally intensive for real-time applications. Using coarser meshes to reduce computation times may prevent the resolution of crucial spatio-temporal scales. Reduced Order Models (ROM) present a compromise between purely data-driven and purely physics-based approaches (see, e.g., [5] for some aeroelastic applications). The Proper Orthogonal Decomposition (POD) -Galerkin is a class of ROM derived from the physical equations but trained using data for the simplified analysis or rapid resolution of particular flow systems, for instance a turbulent mixing layer [6]. The ROM solution is constrained to live inside a small subspace learned from the data. Nevertheless, the unsteady CFD ROM state of the art limits itself to deterministic ROMs (often linear and/or with purely data-driven calibration) with limited prediction capabilities when the ROM dimension is low. It is probably mainly due to the chaotic and intermittent nature of turbulence and closure problems. Indeed, intermittency – ubiquitous in turbulence – is related to rare events and long-memory processes. These long-memory processes make learning from a finite time window more complicated because turbulence data are hardly exhaustive. Therefore, the learned or partially-learned turbulent flow ROMs remain inexact outside the learning time interval and uncontrolled in the long run, owing to their chaotic nature (intrinsic sensitivity to perturbations [7]) and the growth of accumulated error along time. Thus, predictions become less and less accurate. Additionally, we believe that ROM deterministic closures can hardly be accurate in the long run. Indeed, energy fluxes between temporal modes corresponding to orthogonal divergence-free spatial modes (e.g., curl of Fourier modes or POD modes) of a real incompressible flow are described by dyads and triads. The transfer of energy from one temporal mode to another is often dependent on a third temporal mode [8]. It is also true for the transfers of energy from one temporal mode to the mean and from the mean to a temporal mode. Unfortunately, in ROM, the mode truncation breaks many of these triads because the third mode is unknown [9, 10]. The missing negative and positive energy fluxes induce instability and over-damping respectively. In order to stabilize ROMs, either an additional deterministic term (typically an eddy viscosity term) [11, 12, 13] or an additional constraint [7], is often introduced, which may require calibration with the aid of data [14, 15, 16]. However, few authors address the missing positive energy fluxes issue. One reason is that adding relevant terms, which increase energy, is much more difficult in a deterministic framework than in a stochastic one. Nevertheless, these positive energy fluxes are essential to maintain sufficient variability in the linearly-stable temporal modes and thus to maintain coherent ROM dynamics over extended run-time integrations.

The coupling of the model with measurements, namely the data assimilation, presents a number of additional

challenges. Significant advances have been made in the field of meteorology [17]. In the fluid dynamics ROM literature, several variational data assimilation methods have been proposed to jointly learn the ROM (or part of the ROM) [e.g. 18, 14, 19, 20, 15]. However, those variational approaches typically require a large amount of dense data to be assimilated – typically two-dimensional flow observables like particle image velocimetry (PIV). In such a case, fluid flow *estimation* hardly necessitates ROM procedures since the inverse problem is over-determined. A least-square estimation is often sufficient to *estimate* the first velocity temporal modes (i.e. the most energetic modes if POD is considered). This is a kind of velocity mode-based spatial interpolation. For *prediction* of unsteady turbulent flows in a near future, the use of ROM is more relevant, but, as discussed previously, those methods suffer diminished predictive performance outside of the learning interval.

Maintaining satisfactory predictions despite a reduction in the quantity of measurements to assimilate is best realized through better state forecasts or through better adapted data assimilation methods. To obtain better forecasts – i.e. better state prior distributions in an ensemble-based data assimilation framework – improvements to the underlying ROM are required. Data assimilation algorithms fully addressing non-linearities of fluid mechanics are limited by the available computational resources and the dynamical model’s accuracy quantification (uncertainty quantification). The use of ROM with severe dimension reduction alleviates the computational resource limitation but may complicate the uncertainty quantification. Indeed, [21] have already demonstrated that a particle filter can successfully assimilate pointwise measurements into a ROM of wake flows at a Reynolds numbers of 100 and 1000. Nevertheless, this state-of-the-art algorithm suffers from a crude dynamical model’s uncertainty quantification. Accordingly, those good prediction skills require high ROM accuracy and informative measurements, i.e. many modes (8 and 30 modes at a Reynolds numbers of 100 and 1000 respectively) and many measurement points (about the number of modes). Thus, the dimension reduction is not that severe and, as discussed previously, their fluid flow estimation may not require a ROM since many measurements are available. Dynamics under location uncertainty (LU) [22, 23] provide a random fluid mechanics framework designed for improved quantification of dynamic model accuracy. Inspired from the theoretical work of [24], [22] has introduced that stochastic closure and [23] generalized it. It has led to huge improvements in uncertainty and model error quantification both in high-dimensional CFD [25, 26, 27] and reduced state spaces [28, 29]. Nevertheless, no associated ensemble-based data assimilation algorithm has been proposed yet. This paper will be the first on this path. Our new fast data assimilation algorithm for fluid flows includes both an efficient background state ensemble emulator [29] and a particle filter to correct this ensemble.

This paper will be organized as follows: section 2 will recall the main aspects of POD-Galerkin ROM (POD-ROM), section 3 will present the key player of our algorithm: a randomized version of POD-Galerkin ROM, section 4 explains the data assimilation procedure, and finally, section 6 will showcase its potential through some of our numerical results.

## 2. POD-ROM

Reduced Order Models (ROM) aim to reduce the computational expense of simulations by using approximations to significantly diminish the solution's degrees of freedom as compared with the full order model. The approximations are typically achieved through a combination of existing simulation data and modeling based on physical equations. In traditional CFD, the degrees of freedom associated with the velocity fields, for instance, are proportional to the number of grid points in the spatial domain (typically in the order of  $10^6$ ). To achieve a reduction in the *dimensionality* of the solution using ROMs, velocity fields are traditionally decomposed as follows:

$$\mathbf{v}(\mathbf{x}, t) = \underbrace{\mathbf{w}(\mathbf{x}, t)}_{\substack{\text{Resolved} \\ \text{by the ROM}}} + \underbrace{\mathbf{v}'(\mathbf{x})}_{\substack{\text{Unresolved} \\ \text{by the ROM}}}, \quad (1)$$

where the resolved field,  $\mathbf{w}(\mathbf{x}, t)$ , is further partitioned into time-averaged and unsteady components:

$$\mathbf{w}(\mathbf{x}, t) = \underbrace{\bar{\mathbf{v}}(\mathbf{x})}_{\substack{\text{Time} \\ \text{averaged}}} + \underbrace{\sum_{i=1}^n b_i(t)\phi_i(\mathbf{x})}_{\substack{\text{Unsteady} \\ \text{component}}}, \quad (2)$$

with  $1 \leq n \leq 10^2$ . Proper orthogonal decomposition (POD) learns the time-averaged  $\bar{\mathbf{v}}(\mathbf{x})$  and the spatial modes  $\phi(\mathbf{x})$  through principal component analysis (PCA) of a series of high-resolution simulations (training set). Subsequently, physical equations, such as the Navier-Stokes equations, can be projected onto these spatial modes, providing a system of  $n$  coupled ordinary differential equations to describe the evolution of the temporal modes  $b_i(t)$ . Combining temporal integration of this reduced-dimensionality system with a given initial condition and equation (2) enables an *prior* prediction (i.e. a prediction unconstrained by measurements) of the resolved velocity field at any given time. Thus, the ROM calculation scheme represents a compromise between entirely data-driven methods and purely physical models. It combines the available simulation data with physical modeling to achieve reliable predictions with improved efficiency.

## 3. MODEL UNCERTAINTY QUANTIFICATION

Models under location uncertainty (LU) represents a stochastic approach to CFD [22, 25, 26, 27], providing both an efficient ROM *closure* to compensate for neglected degrees of freedom  $\mathbf{v}'$ , and quantification of the errors induced by this closure. LU relies on two assumptions, (1) the time decorrelation of the unresolved velocity component  $\mathbf{v}'$  (see eq. (1)) and (2) the stochastic transport (up to some forcing  $\mathbf{F}$ ) of the resolved velocity component  $\mathbf{w}$ . With Itô stochastic calculus notations (see Appendix A), the Navier-Stokes equation under location uncertainty reads:

$$\frac{Dw_k}{Dt} = \partial_t w_k + \left( \mathbf{w} - \frac{1}{2} (\nabla \cdot \mathbf{a})^T + \mathbf{v}' \right) \cdot \nabla w_k - \frac{1}{2} \nabla \cdot (\mathbf{a} \nabla w_k) = F_k. \quad (3)$$

The unresolved velocity (Eulerian) absolute diffusivity,  $a_{pq}$ , is determined from the unresolved velocity components:

$$a_{pq} = \overline{v'_p v'_q} \tau_{v'}, \quad (4)$$

where  $\overline{v'_p v'_q}$  is the time averaged product of the Eulerian unresolved velocities components of  $v'_p$  and  $v'_q$ .  $\tau_{v'}$  represents the unresolved velocity correlation time. In comparison with classical fluid dynamics conservation equations, models under location uncertainty introduce three additional terms, corresponding to (i) a turbulent diffusion,  $\frac{1}{2} \nabla \cdot (\mathbf{a} \nabla w_k)$ , (ii) a large-scale advecting velocity correction,  $-\frac{1}{2} (\nabla \cdot \mathbf{a})^T$ , and (iii) a multiplicative noise term,  $\mathbf{v}' \cdot \nabla w_k$ . To express the uncertainty induced by the dynamic truncation (inherent to any closure method), multiple simulations can be run in parallel using the stochastic model to efficiently realize the most probable future states of the fluid system [25, 28, 26, 27, 30, 31]. Since this stochastic closure method is based on physics [27], its robustness is proved, and calibrations can be performed from all available physical quantities, including the unresolved velocities  $\mathbf{v}'$ . Almost no tuning nor fitting of the ROM is hence required.

Previous work by the current authors [32] makes use of this formalism in a POD-Galerkin context for data analysis, but without considering the noise term  $\mathbf{v}' \cdot \nabla w_k$ . Here, we do consider this noise term. As in [29], we have implemented the POD-Galerkin of the Navier-Stokes model under location uncertainty (3). We obtain the following ROM:

$$\frac{db_i(t)}{dt} = \mathcal{M}_i \left( \mathbf{b}(t), \dot{\beta}(t) \right) \triangleq \underbrace{\sum_{p=0}^n l_{pi} b_p(t) + \sum_{p=0}^n \sum_{q=0}^n c_{pqi} b_p(t) b_q(t)}_{\text{Usual POD-Galerkin terms}} + \underbrace{\sum_{p=0}^n f_{pi} b_p(t) + \sum_{p=0}^n \sum_{k=1}^n \tilde{\alpha}_{pi k}^R b_p(t) \dot{\beta}_k(t)}_{\text{New POD-LU-Galerkin terms}}, \quad (5)$$

where  $(\dot{\beta}_k)_k$  are  $n$  independent one-dimensional white noises and, by convention, the remaining parameters are  $b_0 = 1$ ,  $\mathcal{M}_0 = 0$  and  $\phi_0 = \bar{v}$ . Using the physical equations (3) and (4) and corresponding technical statistical estimators from stochastic calculus, the ROM coefficients  $l$ ,  $f$ ,  $c$ , and  $\tilde{\alpha}^R$  are determined from the resolved spatial modes  $\phi_i$ , the resolved temporal modes  $b_i$ , and the POD residual velocity  $\mathbf{v}'$ . Interested readers can refer to Appendix B or to [29] for more details. Aimed at applied mathematicians, [29] extensively rely on stochastic calculus notations whereas we have tried to make this paper and its appendices accessible to a broader audience. The inclusion of the noise terms into equation (5), enables the characterization of the model uncertainty and enhances forecasting capabilities beyond the scope of the training window. In the following, we will refer to this ROM as reduced location uncertainty model, abbreviated Red LUM. By extension, our complete data assimilation method will also be referred to as Red LUM.

## 4. PARTICLE FILTERING FOR DATA ASSIMILATION

The final element of Red LUM's pipeline employs a particle filter [33] in order to integrate the real-time measurements from multiple sensors and to synchronize the ROM simulation with the real observed flow. We first explain

115 why particle filtering was preferred to more common approaches, such as Ensemble Kalman Filter (EnKF). We then point out again the principle of this data assimilation algorithm.

The governing equations of fluid mechanics are non-linear and high-dimensional. Solutions are non-Gaussian. Therefore, the Kalman filter, the well-known Gaussian data assimilation technique, does not have good performance. Moreover, the curse of dimensionality prevents the construction of the huge state covariance matrix. Theoretically, 120 the proper method to approach such non-linear and non-Gaussian dynamics is particle filtering [34]. However, this method often requires a large ensemble of forecast realizations, also called particles. When the state dimension increases, the problem quickly becomes intractable as even more particles are needed, and each new particle comes with a severe additional computational cost. For this reason, EnKF or variational data assimilation approaches have historically been preferred, especially for meteorological applications, although recently a number of studies 125 have employed variants of particle filtering to address high-dimensional problems [35, 36, 37, 38, 39]. Variational methods, in particular the widely-used 4D-Var algorithm [40] and its variants, have demonstrated promising performance. However, the 4D-Var algorithm requires adjoint codes and neglects the non-Gaussian and statistical non-stationary nature of the model errors. Furthermore, ensemble-based data assimilation methods are typically easier to parallelize, since particles can be forecast between two assimilation steps independently. The EnKF and 130 its variants (such as the square-root EnKF with localized and inflated covariance) are less sensitive to the curse of dimensionality than particle filters, but rely on linear state corrections, which can lead to non-physical solutions [27]. Here, a fully-nonlinear filter – the particle filter – can be used for two reasons. Firstly, the POD-ROM makes the state dimension small enough to alleviate the curse of dimensionality. Secondly, models under location uncertainty efficiently spread small ensembles of particles over the state space without introducing new errors [28, 27, 29]. 135 Therefore, particle filtering accommodates a small ensemble in this context, which greatly reduces its computational cost.

Algorithm 1 presents an overview of the commonly used sequential importance resampling (SIR) particle filter [41, 42] employed in this study. The first step initializes an ensemble of  $N_p$  independent states which are assigned equal weights. The independent states are referred to as particles or sometimes realizations. Particles are forecast 140 in time using the evolution model (5). On each assimilation of a new measurement, the weights are updated based on the particles' likelihoods. An additional *re-sampling* step prevents the weights variance from increasing over time (degeneracy or particle impoverishment), which results in poor state estimations. Ultimately, the ensemble forms the shape of the posterior distribution (i.e. the state distribution conditioned on the assimilated data) and gives the state estimation. The particles' likelihoods computation depends on the type of assimilated measurements. The 145 next section deals with these computations for the measurements considered in this study.

---

**Algorithm 1** Particle Filter SIR with Red LUM
 

---

**Initialization**

- Compute the ROM coefficients  $l$ ,  $f$ ,  $c$ , and  $\tilde{\alpha}^R$  from a simulation output dataset ▷ POD-Galerkin  
(see [Appendix B](#) for definitions).
- Compute matrices  $\mathbf{A}$  and  $\mathbf{B}$  (see [Appendix D](#) for definitions). ▷ Log-likelihood matrices
- Sample  $\begin{pmatrix} b_1^{(j)}(0) \\ \vdots \\ b_n^{(j)}(0) \end{pmatrix} \stackrel{iid}{\sim} \mathcal{N} \left( 0, \begin{pmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_n \end{pmatrix} \right)$  where  $\lambda_i = \overline{b_i^2}$ . ▷ Initializing the first state

**Loop over time  $t$** 
*Importance sampling*

- $\mathbf{b}^{(j)}(t) = \mathbf{b}^{(j)}(t - dt) + \mathcal{M} \left( \mathbf{b}^{(j)}(t - dt), \dot{\mathbf{b}}^{(j)}(t - dt) \right) dt$ . ▷ State transition
- If an observation  $y(t)$  is available at the current time  $t$ :
  - $l_j(t) = l \left( \mathbf{y}(t) | \mathbf{b}^{(j)}(t) \right) = (\mathbf{b}^{(j)}(t))^T \mathbf{A} \mathbf{b}^{(j)}(t) + \mathbf{y}(t)^T \mathbf{B} \mathbf{b}^{(j)}(t)$ ; ▷ Log-likelihood up to a constant  
(see [Appendix D](#) for the proof of this formula)
  - $l_j(t) = l_j(t) - \max_j l_j(t) + 90$ ; ▷ Add a constant to prevent numerical errors when applying exp
  - $W_j(t) = \exp(l_j(t))$ ; ▷ Computing weights
  - $\mathbf{W}_j(t) = \frac{W_j(t)}{\sum_{m=1}^{N_p} W_m(t)}$ ; ▷ Normalization

*Re-sampling*

- Each new temporal mode  $b^{(j)}(t)$  is replaced by one of the old temporal modes ▷ Resampling  
 $b^{(1)}(t), \dots, b^{(N_p)}(t)$  with probability  $\mathbf{W}_1(t), \dots, \mathbf{W}_{N_p}(t)$ , respectively.

**Final posterior distribution** at a time  $t$  larger than measurement times  $t_1, \dots, t_K$

$$p(\mathbf{b}(t) | \mathbf{y}(t_1), \dots, \mathbf{y}(t_K)) \approx \sum_{j=1}^{N_p} \frac{1}{N_p} \delta \left( \mathbf{b}(t) - \mathbf{b}^{(j)}(t) \right). \quad \text{▷ Posterior Distribution}$$


---

## 5. MEASUREMENTS TO ASSIMILATE

In theory, Red LUM can assimilate any measurements. Here, we choose a widely-used fluid flow velocimetry technique: particle image velocimetry (PIV). The measurement method uses a high-power light source and high-speed camera to acquire images of solid tracer particles immersed in the fluid. A combination of image processing and optical flow algorithms are employed to characterize the spatio-temporal velocity distribution of the fluid.

Specifically, our algorithm was tested using a cropped two-dimensional, two-component PIV (2D2C PIV) experimental configuration. In order to assimilate measurements, we need to mathematically model the link between the system state – the full three-dimensional three-component velocity  $\mathbf{v}$  or here its reduced representation  $\mathbf{b}$  – and the measurements – here the PIV field. Such a mathematical model is called observation model and is used to compute the particle filter log-likelihoods  $l_j(t)$  (see algorithm 1 and Appendix D). We here propose the following linear observation model:

$$\mathbf{y} = \mathcal{H}[\mathbf{v}] + \boldsymbol{\epsilon}_y, \quad (6)$$

$$= \sum_{i=0}^n \mathcal{H}[\phi_i] b_i + \underbrace{(\mathcal{H}[\mathbf{v}'] + \boldsymbol{\epsilon}_y)}_{=\boldsymbol{\epsilon}_y^R(t)}, \quad (7)$$

where  $\mathbf{y}$  is the raw PIV and  $\boldsymbol{\epsilon}_y$  represents the PIV measurement error. The linear operator  $\mathcal{H}$  incorporates a 3-dimensional spatial smoothing operation as well as occlusion of the horizontal plane, and its corresponding component in the velocity field, to approximate the PIV measurements. The parameters inside  $\mathcal{H}$  and  $\boldsymbol{\epsilon}_y$  are estimated using experimental data, comparing the hot-wire and PIV measurements' spectrum (using a Taylor assumption). Appendix C details the definition of the  $\mathcal{H}$  operator and the PIV measurement noise covariance.

Additionally, to make the data assimilation task more challenging, the information relating to a large subset of points in the grid was obscured through the operator  $\mathcal{H}$ . This results in a small observation vector  $\mathbf{y}$ . Indeed, estimating a vector  $b$  of  $n \sim 10$  components from a noisy linearly-dependent observation vector  $\mathbf{y}$  of  $M_{PIV} \sim 10^4 (\gg n)$  components is often an over-determined inverse problem and could otherwise be solved using a straightforward least-squares procedure.

Note that the strong influence of the unresolved velocity  $\mathbf{v}'$  on the final observation model's uncertainty  $\boldsymbol{\epsilon}_y^R(t)$  is taken into account through a noise term in (7). This will naturally influences the data assimilation algorithm results, since the particle filter log-likelihoods  $l_j(t)$  are computed from this reduced observation model (see algorithm 1 and Appendix D).

This study considers synthetic PIV data, generated applying the general observation model (6) to fully-resolved CFD simulation outputs  $\mathbf{v}$ . The fully-resolved simulations provide a complete instantaneous 3D characterization of the velocity field for the data assimilation validation. Note these synthetic PIV measurements are not generated by the reduced observation model (7), which is used for the particle filter algorithm only.

## 6. NUMERICAL RESULTS

A performance assessment of our reduced data assimilation algorithm (Red LUM) will be conducted in relation to two distinct cylindrical wake flows. The two wake flows correspond to two- and three-dimensional flows past an isolated circular cylinder at Reynolds numbers of 100 and 300 respectively. Cylindrical wake flows exhibit pseudo-periodic vortex shedding cycles at the rear surface of the cylinder. This complex flow behavior provides an excellent basis to assess alternative flow prediction methods. In addition to Red LUM, data assimilation of the wake flows measurements will also be presented for two state-of-the-art POD-ROMs, detailed in section 6.1, to provide a basis of comparison. The predictive performance of each model will first be reviewed through a qualitative assessment of the vorticity and the Q-criterion of the reconstructed velocity fields, prior to a quantitative assessment of the errors associated with the reconstructed velocity fields. Subsequently, sensitivity analysis will be performed to determine the influence of the ensemble size and the size and location of the assimilated data on the model performance. Finally, we provide some insights into our algorithm complexity.

Prior to investigation of the various POD-ROM methods, reference simulations were performed using direct numerical simulations (DNS), implemented using Incompact3d, a high-order flow solver based on the discretization of the incompressible Navier-Stokes equations [43]. These fully-resolved simulations are illustrated in top of figure 1 and top right of figure 3. The two wake flows exhibit major differences in their complexity and number of degrees of freedom. At a Reynolds number of 100, the flow remains laminar, whereas, at a Reynolds number of 300, the flow appears to be three-dimensional and much more complex. The full-order simulation's spatial grids at Reynolds numbers of 100 and 300 correspond to state-space dimensions of about  $10^4$  and  $10^7$ , respectively. During the learning period, 140 and 80 vortex shedding cycles are used to construct the ROMs. The remaining vortex shedding cycles are used to build synthetic measurements and to validate the method. For both flows, a single spatial resolution point of the synthetic PIV data is assimilated ten times for each vortex shedding cycle. The observation point, indicated as a star in figure 1, is located just above the recirculation zone, at coordinates of  $x = 1.31D$  (streamwise direction),  $z = 0$  (spanwise direction), and  $y = 1.27D$  (orthogonal direction). The cylinder, of diameter  $D$ , is centered at  $(0, 0, -)$ . We refer to this setting as observation case 1. Other observation cases will be considered at the end of this section.

### 6.1 State-of-the-art ROMs

The Red LUM performance will be assessed against two alternative state of the art algorithms. To facilitate a fair comparison, the same observation model (7), assimilation method (SIR particle filter), and ensemble size (100 particles) will be employed in each case.



### 6.1.1 Deterministic state-of-the-art ROM (D-SOTA)

The first state-of-the-art ROM is a usual deterministic POD-Galerkin model, with an optimally fitted eddy (using a least square method), presented in (8):

$$\frac{db_i(t)}{dt} = \underbrace{\sum_{p=0}^n l_{pi} b_p(t) + \sum_{p=0}^n \sum_{q=0}^n c_{pqi} b_p(t) b_q(t)}_{\text{Usual POD-Galerkin terms}} + \underbrace{\left( \frac{\nu^{\text{ev}}}{\nu} - 1 \right) \sum_{p=0}^n l_{pi} b_p(t)}_{\text{Fitted on the resolved temporal modes dynamics}}, \quad (8)$$

where  $\nu^{\text{ev}}$  is the fitted eddy viscosity, and  $\nu$  is the molecular viscosity. While this POD-ROM method is itself deterministic, ensemble-based data assimilation can nonetheless be implemented through randomized initial conditions.

### 6.1.2 Stochastic state-of-the-art ROM (S-SOTA)

The randomization of initial conditions has been shown to typically underpredict errors in fluids dynamics forecasts [44, 45, 46, 47] and so an alternative stochastic state-of-the-art ROM is investigated for comparison. The second baseline POD-ROM adds a white noise term to the first baseline ROM. The addition of white Gaussian noise forcing is a simple *ad hoc* randomization technique for a given deterministic dynamical system [e.g. 48, 49]. Despite its potential lack of physical relevance, such a strategy is very often adopted in data-assimilation applications [28].

$$\frac{db_i(t)}{dt} = \underbrace{\sum_{p=0}^n l_{pi} b_p(t) + \sum_{p=0}^n \sum_{q=0}^n c_{pqi} b_p(t) b_q(t)}_{\text{Usual POD-Galerkin terms}} + \underbrace{\left( \frac{\nu^{\text{ev}}}{\nu} - 1 \right) \sum_{p=0}^n l_{pi} b_p(t) + \sum_{k=1}^n \sigma_{ik}^{\text{ev}} \dot{\beta}_k(t)}_{\text{Fitted on the resolved temporal modes dynamics}}, \quad (9)$$

where  $\nu^{\text{ev}}$  is the fitted eddy viscosity,  $\nu$  is the molecular viscosity, and  $(\sigma^{\text{ev}})(\sigma^{\text{ev}})^T$  is the fitted additive noise covariance matrix. For more details about this state-of-the-art stochastic ROM, one can refer to [29] for example.

[21] also proposed a similar POD-ROM for application with a particle filter. During the training phase, their algorithm estimates a whole additional linear term rather than just an eddy viscosity coefficient. Moreover, their noise term forces the POD-ROM at each assimilation step only (as opposed to every simulation time step) and employs a tuned variance. The (reduced) observation model benefits from greater simplicity as the observation noise  $\epsilon_y^R(t)$  has a tuned spatially-uniform variance.

## 6.2 Qualitative performance analysis

A performance assessment of each POM-ROM method will be conducted using the DNS predicted velocity field, denoted  $\mathbf{v}_{\text{dns}}$ , as a full-order reference case. However, the solution of each POD-ROM method is confined to the affine space spanned by the POD modes (a set of fields defined according to (2)). Consequently, the optimum<sup>1</sup> prediction which can be achieved using a POD-ROM-based method (or theoretical performance limit) corresponds

---

<sup>1</sup>Best in terms of  $L^2$  norm. Note that by definition, the POD is also the best modal decomposition in terms of  $L^2$  norm for a given degree of freedom  $n$ .

to the orthogonal projection of the DNS field onto this reduced affine space. This theoretically optimum field,  $\tilde{\Pi}_\phi[\mathbf{v}_{\text{dns}}]$ , is considered in the following as the reference field. It is defined as follows:

$$\tilde{\Pi}_\phi[\mathbf{v}_{\text{dns}}] \triangleq \bar{\mathbf{v}} + \sum_{i=1}^n b_i^{\text{ref}} \phi_i \quad \text{with} \quad b_i^{\text{ref}} \triangleq (\phi_i, \mathbf{v}_{\text{dns}} - \bar{\mathbf{v}}). \quad (10)$$

This projected field corresponds to the known temporal modes  $b_i(t) = b_i^{\text{ref}}(t)$  ( $i \leq n$ ) in the equation (2). These projected fields are indicated in the second panel of figures 1 and 2 for a Reynolds number of 100 and at the top left of figures 3 and 4 for a Reynolds number of 300 respectively. Note that when the Reynolds number increases, even though the projected velocity field is still meaningful, it can be relatively far from the DNS field.

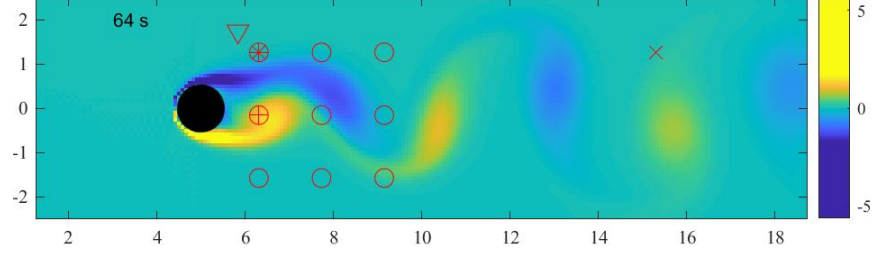
### 6.2.1 Reynolds number of 100

Figures 1 and 2 present estimated vorticity fields<sup>2</sup> of the data assimilation results at a Reynolds number of 100 and  $t = 64 s$  (after the training period), with  $n = 2$  and 8 modes, respectively. The four simulation results presented correspond to the 2D DNS, the DNS projection onto the POD modes (corresponding to the theoretical performance limit of POD-ROM simulations), Red LUM, and S-SOTA (the POD-Galerkin with optimally fitted eddy viscosity and additive noise, see (9)). From an initial qualitative assessment of the flow field, Red LUM estimations appears to closely resemble the DNS results and demonstrates excellent agreement with the theoretical performance limit, indicating excellent predictive potential beyond the learning window in spite of the limited assimilated data (single measurement point). Quadrupling the number of modes from 2 to 8 does not appear qualitatively to compromise the predictive performance of the proposed POD-ROM. In order to better appreciate the Red LUM potential, figures 1 and 2 also display the predictions from the stochastic state-of-the-art POD-ROM (S-SOTA) (9) (the POD-Galerkin with optimally fitted eddy viscosity and additive noise). The fitted additive noise variance of this POD-ROM is necessarily large in order to encompass the significant ROM error. Consequently, the prior probability distribution (i.e. the probability distribution of the (background) state before taking account the measurement) generated by S-SOTA is hardly informative, while the single measurement gives little additional information. These limitations result in a departure from the DNS predictions which becomes particularly apparent as the order of the ROM is increased from 2 to 8 and the fitted additive noise variance is enhanced to accommodate the increased complexity of the flow dynamics. The results of the deterministic POD-ROM method (D-SOTA) (8) were less promising than the stochastic approach and so the vorticity fields are not presented. Instead, a quantitative assessment of D-SOTA will be presented later in section 6.3. In contrast, the physical structure of Red LUM, and in particular of its skew-symmetric multiplicative noise, guarantees an efficient prior probability distribution [29], enabling more accurate data assimilation.

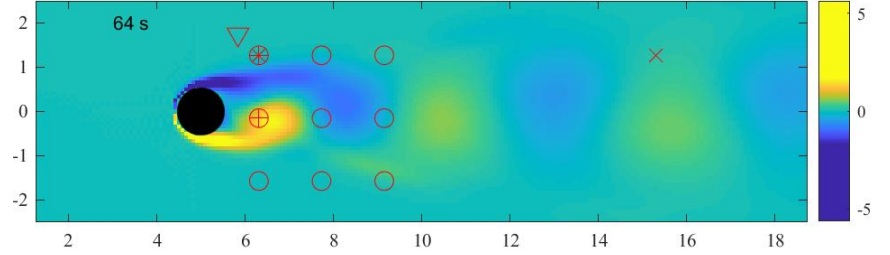
---

<sup>2</sup>The vorticity field is the curl of the velocity field commonly used to visualize 2D vortices.

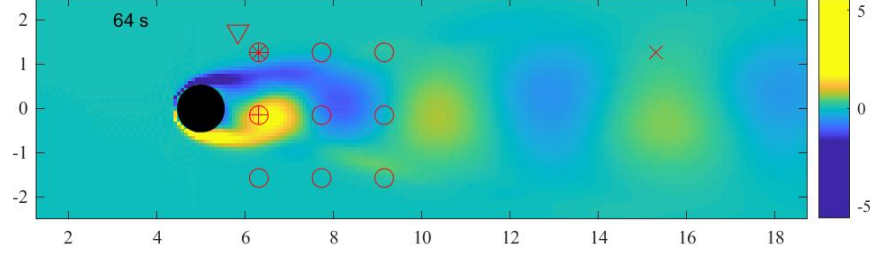
Reference full-order simulation  $\mathbf{v}_{\text{dns}}$   
 (2D DNS at Reynolds 100:  
 state space dimension  
 of about  $10^4$ )



Optimum prediction:  
 Projection of the DNS  
 onto the POD basis  $\tilde{\Pi}_\phi[\mathbf{v}_{\text{dns}}]$   
 (State space  
 of dimension 2)



Red LUM:  
 Our reduced data assimilation prediction  
 (ROM state space  
 of dimension 2)



S-SOTA:  
 State-of-the-art method prediction  
 (ROM state space  
 of dimension 2)

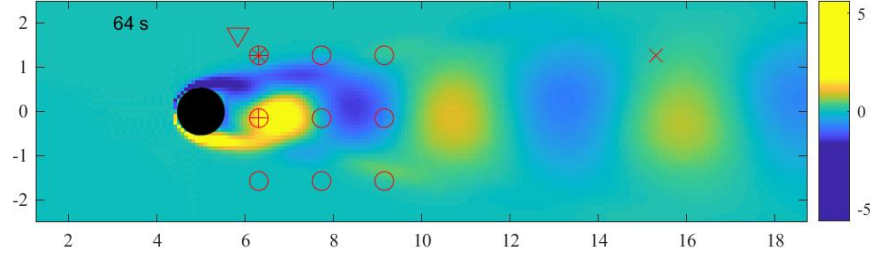
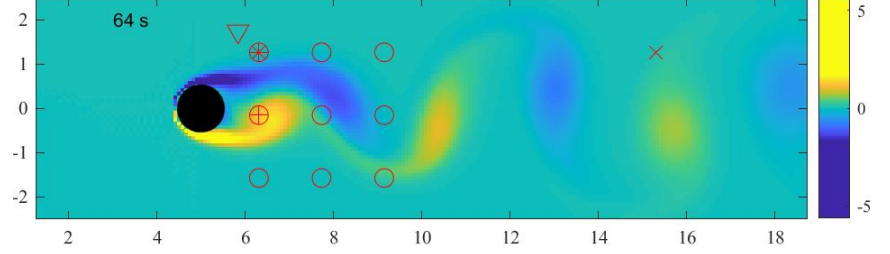
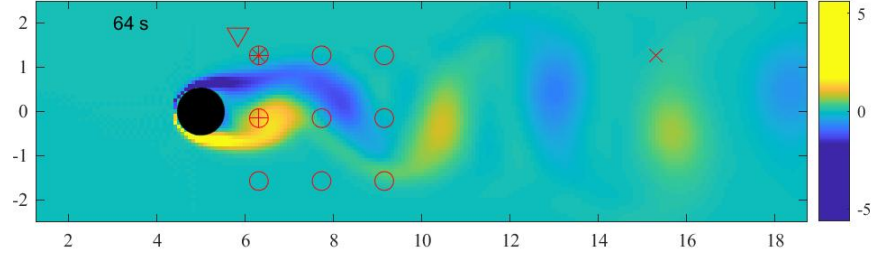


Figure 1: Vorticity field – of (from top to bottom) the 2D DNS at a Reynolds number of 100, the reference 2-dimensional representation (projection of the DNS onto the POD modes), Red LUM, and S-SOTA. The red signs indicate the measurement locations for the different observation cases considered: case 1 (star symbol) (considered in sections 6.2 and 6.3 and figures 1-5 and 8), case 2 (circle symbols), case 3 (cross symbol close to the cylinder) (both considered in section 6.4.1 and figure 6), case 4 (triangle symbol), and case 5 (downstream cross symbol) (both considered in section 6.4.1 and figure 7).

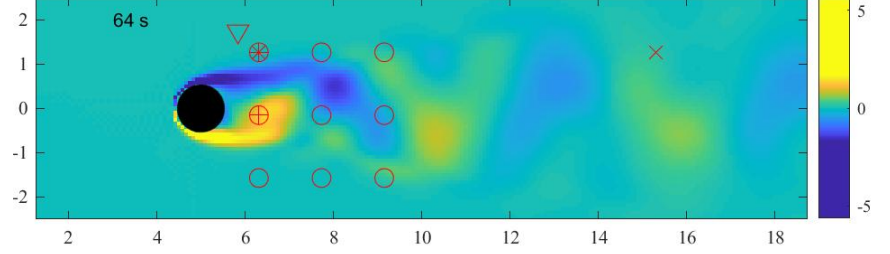
Reference full-order simulation  $\mathbf{v}_{\text{dns}}$   
 (2D DNS at Reynolds 100:  
 state space dimension  
 of about  $10^4$ )



Optimum prediction:  
 Projection of the DNS  
 onto the POD basis  $\tilde{\Pi}_\phi[\mathbf{v}_{\text{dns}}]$   
 (State space  
 of dimension 8)



Red LUM:  
 Our reduced data assimilation prediction  
 (ROM state space  
 of dimension 8)



S-SOTA:  
 State-of-the-art method prediction  
 (ROM state space  
 of dimension 8)

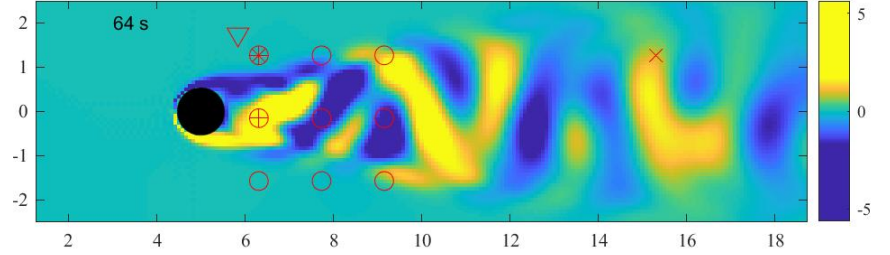


Figure 2: Vorticity field – 13 vortex shedding cycles after the learning period – of (from top to bottom) the 2D DNS at a Reynolds number of 100, the reference 8-dimensional representation (projection of the DNS onto the POD modes), Red LUM, and S-SOTA. The red signs indicate the measurement locations for the different observation cases considered: case 1 (star symbol) (considered in sections 6.2 and 6.3 and figures 1-5 and 8), case 2 (circle symbols), case 3 (cross symbol close to the cylinder) (both considered in section 6.4.1 and figure 6), case 4 (triangle symbol), and case 5 (downstream cross symbol) (both considered in section 6.4.1 and figure 7).

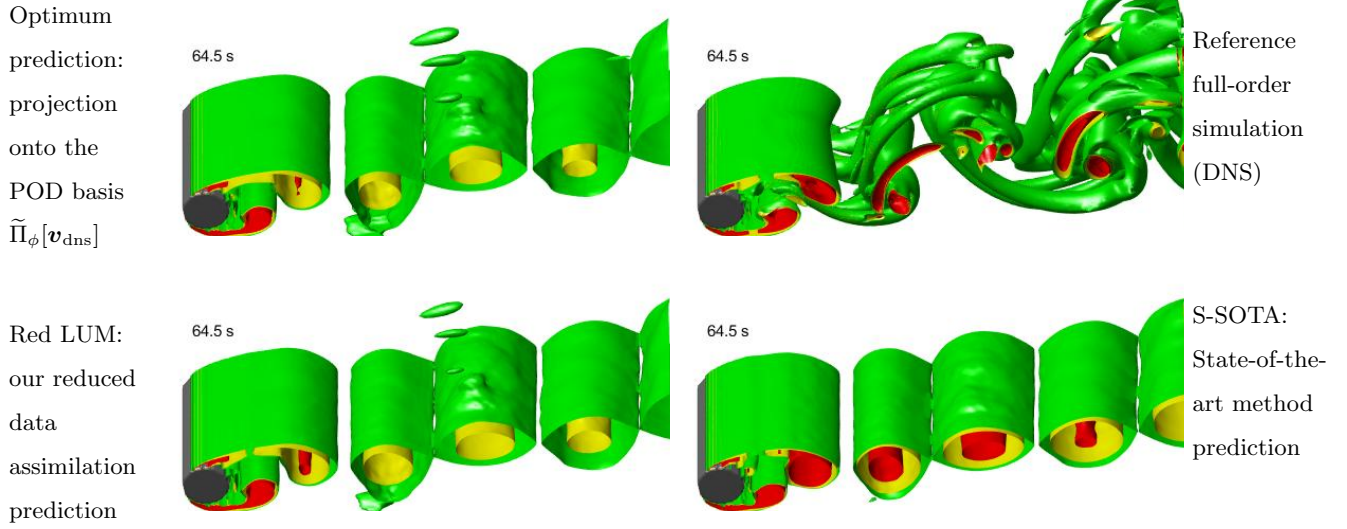


Figure 3: Q-criterion – 13 vortex shedding cycles after the learning period – from the reference 2-dimensional representation (projection of the 3D DNS at Reynolds 300 onto the POD modes) (top left), the DNS (top right), Red LUM (bottom left), and S-SOTA (bottom right).

### 6.2.2 Reynolds number of 300

Figures 3 and 4 show the Q-criterion<sup>3</sup> isosurfaces of the data assimilation results at a Reynolds number of 300, with  $n = 2$  and 8 modes, respectively. As with the 2D laminar case, the proposed estimations closely resemble the theoretical POD-ROM performance limits. These projected references (top left plots in figures 3 and 4) correspond to the known temporal modes  $b_i(t)$  ( $i \leq n$ ) in the equation (2). The qualitative departure between these optimums estimations and the DNS reference correspond to the unresolved velocity  $v'$ . This velocity component is mainly restricted to small-scale 3-dimensional effects at a Reynolds number of 300. S-SOTA demonstrates greater departure from both the DNS and the theoretical performance limit since, as for the Reynolds number of 100, the noise structure of S-SOTA is not well adapted to the problem under investigation.

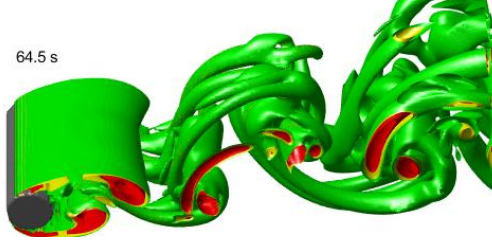
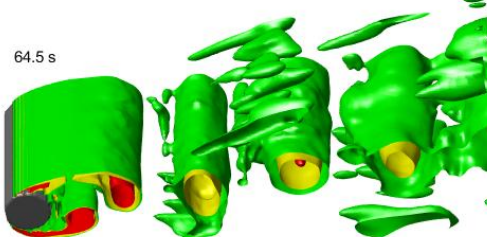
## 6.3 Quantitative performance analysis

Vorticity fields and Q-criterion isosurfaces provide useful qualitative assessments of the estimated flow fields, however quantitative analysis is required for a comprehensive understanding of the model performance and limitations. The POD-ROM performance will be assessed using a global velocity estimation normalized error. The estimated velocity field is denoted  $\mathbf{w}^{\text{est}} = \sum_{i=0}^n b_i^{\text{est}} \phi_i$ . The Mean Square Error (MSE) simplifies as follows, due to the orthogonality

<sup>3</sup>Q-criterion is a tool used for the visualization of vortices in 3D CFD.

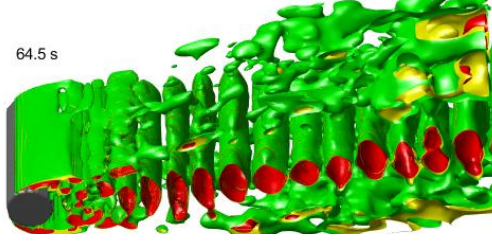
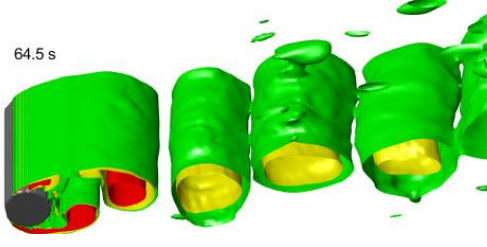


Optimum  
prediction:  
projection  
onto the  
POD basis  
 $\tilde{\Pi}_\phi[\mathbf{v}_{\text{dns}}]$



Reference  
full-order  
simulation  
(DNS)

Red LUM:  
our reduced  
data  
assimilation  
prediction



S-SOTA:  
State-of-the-  
art method  
prediction

Figure 4: Q-criterion – 13 vortex shedding cycles after the learning period – from the reference 8-dimensional representation (projection of the 3D DNS at Reynolds 300 onto the POD modes) (top left), the DNS (top right), Red LUM (bottom left), and S-SOTA (bottom right).

properties of the orthogonal projection  $\tilde{\Pi}_\phi[\mathbf{v}_{\text{ref}}]$  and of the POD modes  $\phi_i$ :

$$\text{MSE} \triangleq \int_{\Omega} \|\mathbf{w}^{\text{est}} - \mathbf{v}_{\text{dns}}\|^2, \quad (11)$$

$$= \int_{\Omega} \underbrace{\|\mathbf{w}^{\text{est}} - \tilde{\Pi}_\phi[\mathbf{v}_{\text{dns}}]\|^2}_{=\sum_{i=1}^n (b_i - b_i^{\text{ref}})^2} + \int_{\Omega} \underbrace{\|\mathbf{v}_{\text{dns}} - \tilde{\Pi}_\phi[\mathbf{v}_{\text{dns}}]\|^2}_{=\mathbf{v}'}, \quad (12)$$

$$= \underbrace{\sum_{i=1}^n (b_i - b_i^{\text{ref}})^2}_{\text{ROM-dependent}} + \underbrace{\int_{\Omega} \|\mathbf{v}_{\text{dns}} - \tilde{\Pi}_\phi[\mathbf{v}_{\text{dns}}]\|^2}_{\text{ROM-independent}}, \quad (13)$$

where the integration is performed over the spatial domain denoted  $\Omega$  and  $\|\bullet\|$  represents the usual Euclidean norm of  $\mathbb{R}^2$  or  $\mathbb{R}^3$ . The first error term will be dependent on the ROM and associated data assimilation method, whereas the second error term depends solely on the POD modes  $\phi_i$ . Consequently, the second error term is the same for Red LUM, D-SOTA and S-SOTA. Note that this term is time-dependent in general. Here, the MSE is normalized by the mean kinetic energy (MKE) in the far fluid moving frame (averaged over the training set),  $\text{MKE} = \overline{\int_{\Omega} \|\mathbf{v} - \mathbf{v}_{\infty}\|^2}$ , where  $\mathbf{v}_{\infty} = (1 \text{ m.s}^{-1}, 0, 0)^T$  is the velocity far from the cylinder. Subtracting the far fluid frame velocity,  $\mathbf{v}_{\infty}$ , renders the MKE, and thus the proposed RMSE (its square root), independent of the spatial domain  $\Omega$  for a sufficiently large spatial domain. It makes our normalized RMSE also independent of the spatial domain choice (for a sufficiently large spatial domain) and hence makes it more objective. However, it should be noted that  $\overline{\int_{\Omega} \|\mathbf{v}\|^2}$ , even normalized by the volume or the surface, remains dependent on the spatial domain since the velocity

variability (the wake) is highly localized.

Figure 5 presents the global velocity estimation normalized error at Reynolds numbers of 100 and 300, between the conclusion of the training period and the termination of the 20<sup>th</sup> and 14<sup>th</sup> vortex shedding cycle respectively. The orange, light blue and dark blue profiles correspond to the errors associated with Red LUM, S-SOTA, and D-SOTA respectively. *Posterior* standard deviations<sup>4</sup> are also visible in shaded colors, and quantify the ensembles diversities after data assimilation. Accordingly, they should be proxies for the velocity estimations error amplitudes. S-SOTA errors appear very large and regularly exceed the black profile corresponding the condition that all temporal modes are set to zero. Furthermore, the S-SOTA errors escalate as the number of modes is increased, whereas the D-SOTA performance quickly diminishes over time (with increased departure from the learning window), particularly at a Reynolds number 300. Conversely, Red LUM error appears relatively stable over time and the magnitude of the errors remains modest even as the number of modes is increased to 8. D-SOTA limitations are explained by the underestimation of the variance outside of the learning window leading to degeneration of the filter, a prediction with escalating errors with increased departure from the learning window, and – at a Reynolds number of 300 – divergence in time of the POD-ROM estimation (not shown). This variance decay illustrates the classical problem of the alignment of ensembles along unstable directions [50, 51], which results from the over-damping of stable temporal modes [9, 10]. This over-damping and, more generally, the variance decay is induced by the missing positive energy fluxes toward each temporal mode (the mode truncation removes many triad-based fluxes) and the stabilizing corrective terms (the eddy viscosity term in this instance).

## 6.4 Influence of the available resources

In practice, the result of our algorithm can be influenced by several factors, such as the quantity and the quality of the measurements available for assimilation, or limitations in the available computational resources.

### 6.4.1 Influence of the available measurements

The dependence of model performance on the data available for assimilation will be assessed in relation to the more complicated 3D wake flow at a Reynolds number of 300. The various observation points used for this analysis are detailed in table 1. Figures 6 and 7 present the corresponding estimation errors for S-SOTA and Red LUM.

Predictably, the performance of each of the methods improves when more measurements are available for assimilation (see observation case 2, left panel of figure 6). S-SOTA remains less efficient than Red LUM, although the predictions tend to converge as the number of observation points is increased, typically as the number of observations  $M_y$  exceeds the dimension  $n$  of the reduced space. This conclusion is consistent with the findings of [21]

---

<sup>4</sup>the square roots of the global velocity estimations variances conditioned on the assimilated data, integrated over the spatial domain, and adequately normalized

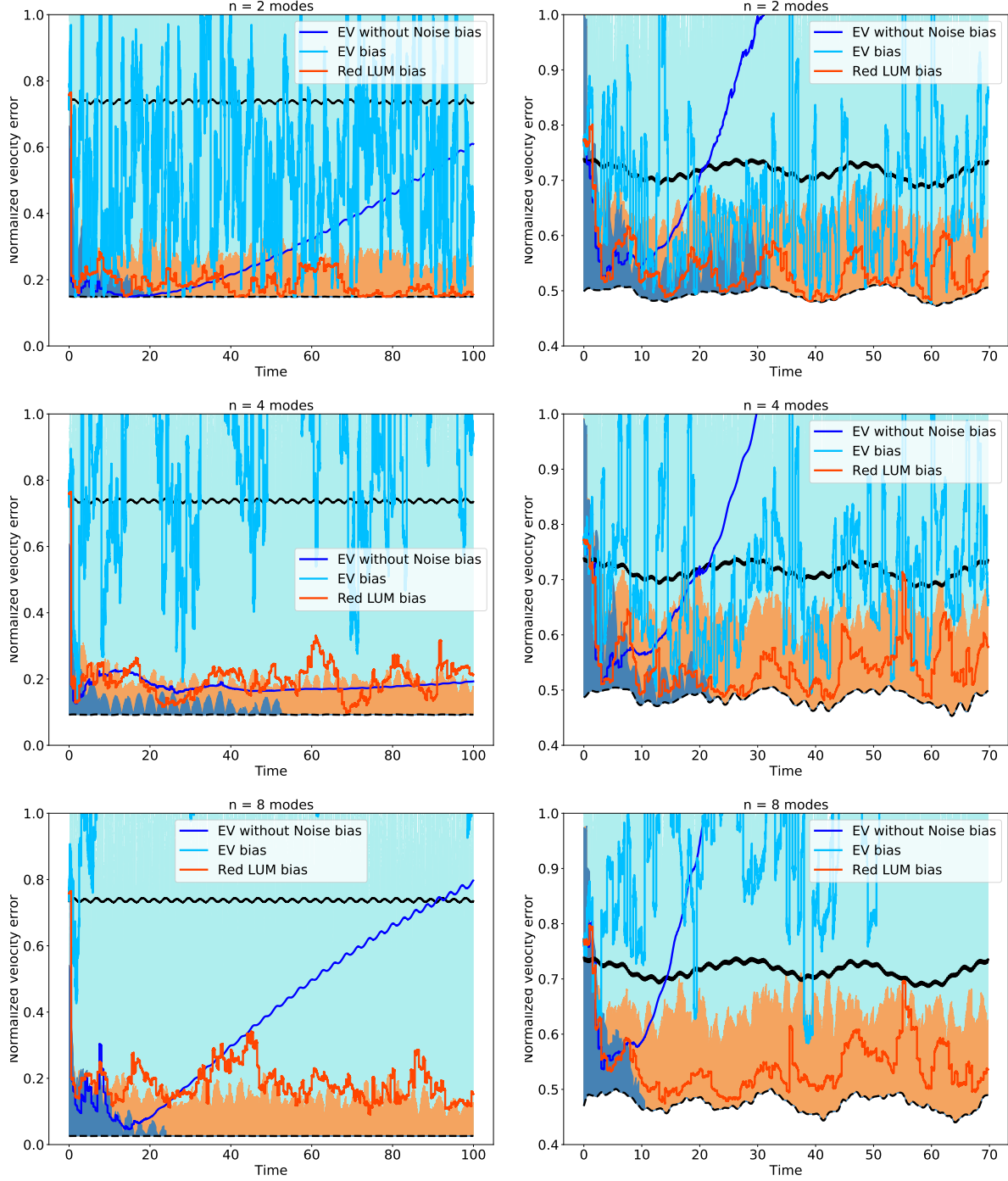


Figure 5: Global velocity prediction normalized error (after the learning period) for the wake flow at Reynolds 100 (left) and 300 (right) in the observation case 1 with, from top to bottom  $n = 2, 4$  and 8 modes for Red LUM (orange), S-SOTA (light blue), and D-SOTA (dark blue). The shaded colors (light orange, light blue, and grey) correspond to the respective estimated *posterior* standard deviations. The dashed black line at the bottom is the POD truncation error. The solid black line at the top is the error obtained by setting all temporal modes  $b_i$  to 0, i.e. keeping only the time averaged velocity  $\bar{v}$ .



Case name	Description	$x/D$	$y/D$	$z/D$	Symbol on figures 1 and 2
Case 1	1 observation point near the recirculation zone	1.31	1.27	0	*
Case 2	$3 \times 3$ observation points	1.31	1.27	0	o
		2.73	1.27	0	
		4.15	1.27	0	
		1.31	-0.15	0	
		2.73	-0.15	0	
		4.15	-0.15	0	
		1.31	-1.57	0	
		2.73	-1.57	0	
		4.15	-1.57	0	
Case 3	1 observation point inside the recirculation zone	1.31	-0.15	0	+
Case 4	1 point outside the wake	0.84	1.74	0	$\nabla$
Case 5	1 observation point far downstream	10.31	1.27	0	$\times$

Table 1: Details of the observation cases considered for the assessment of the methods dependence on the quantity of measurements to assimilate.

which demonstrated excellent flow predictions using a particle filter in conjunction with a POD-ROM similar to S-SOTA employing 9 observations points. However, when the number of observation points is excessively large, the standalone observation model represents an over-determined least-square problem (the number of equations exceeds the number of unknown variables), and so the requirement for a dynamic ROM for flow estimation is less apparent. Note also that when the number of observation points is greatly increased, the likelihood becomes a very peaky function of the state (i.e. the likelihood support is very localized). Accordingly, the particle filtering algorithm needs to be adapted to prevent filter degeneracy. For the wake flow under consideration, tempering and non-Gaussian jittering [35, 38] have been shown to provide efficient solutions at around  $M_y \sim 10^4$  observations points (not shown). Considering a single observation point may be considered as an over-complicated estimation problem compared to practical applications. However, real industrial applications are often required to assimilate data from limited, and sometimes low quality, sensors which provide relatively poor information on the flow characterization. Additionally, many parameters (e.g., boundary conditions, Reynolds number) are either unknown or poorly characterized for many realistic situations and must be estimated on top of the state variables. Consequently, the available information is often insufficient to constraint the estimation problem. For this reason, the extent to which the requirement for external data can be minimized, perhaps to a single measurement as the proposed model aims to achieve, though the incorporation of more physics into the POD-ROM and associated ensemble forecasts, is potentially highly valuable for industrial applications.

The performances of the estimation algorithms are not only sensitive to the quantity of data available for assimilation but the location of the observation points. For an observation point located within the recirculation zone (see observation case 3 presented in the right panel of figure 6), S-SOTA performance marginally improves however the RMSE remains close or above the zero solution RMSE indicates a comparable performance to the mean-velocity solution, whereby all temporal modes  $b_i$  tend to 0.

The Red LUM method exhibits similar flow predictions using observation points 1 and 3, and thus appears less sensitive to the precise location of the chosen observation point. For an observation point outside the wake (observation case 4, left panel of figure 7), all method skills strongly deteriorate, as expected. For an observation point far downstream (observation case 5, right panel of figure 7), S-SOTA keeps similar results. The skills of Red LUM slightly deteriorate but remain much better than S-SOTA. To generalize, the flow predictions appear satisfactory so long as the observation point contacts the wake but most accurate for observations close to the recirculation zone.

Additional tests were performed using assimilation data acquired at a diminished frequency, specifically a solitary measurement per vortex shedding cycle. As with a reduction in the number observation points, reducing the temporal frequency of data acquisition diminishes the performance of the flow estimation. For Red LUM, assimilating data from  $3 \times 3$  observation points at a rate of one per vortex shedding cycle exhibits a comparable performance to simulations based on a solitary observation point with five times the data acquisition rate.

#### 6.4.2 Influence of the number of realizations

Figure 8 presents the dependence of the RMSE on the ensemble size  $N_p$  for Red LUM and for S-SOTA at  $n = 2, 4$  and 8 modes. As the ROM dimension  $n$  is increased, larger ensemble sizes are typically required to achieve an acceptable model error. In each case, no more than 100 realizations are required for Red LUM to converge to its optimum prediction. The rapid convergence of Red LUM in comparison with S-SOTA provides another indication that the structure of the noise is better adapted to the fluid dynamics, and thus fewer particles are required to investigate the most probable flow predictions. The requirement for a smaller ensemble size to converge to a stable error of acceptable magnitude, can incur a significant saving in computational power and thus demonstrates a greatly enhanced potential for real-time applications.

### 6.5 Discussion on the implementation and run-time considerations

Several hours of simulation using on a supercomputer are necessary to generate the training simulation dataset with the off-line high-resolution highly-optimized CFD code.

By contrast, the off-line ROM construction is performed in a few hours on a laptop using the original non-parallelized MATLAB code employed during this study. The POD-ROM building algorithm has now been implemented in C++ using the OpenFOAM-based library ITHACA-FV [52]. While the C++ code remains non-

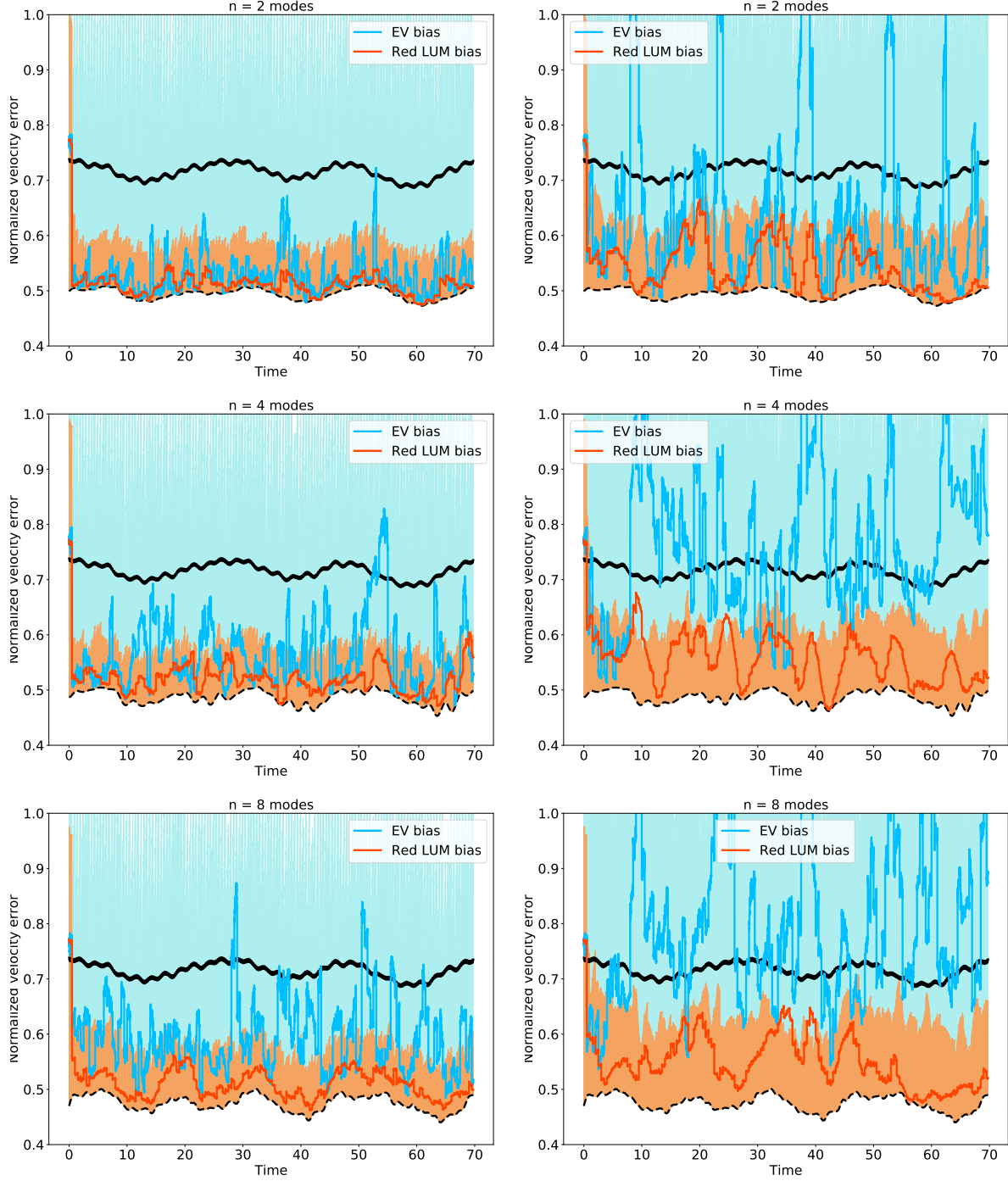


Figure 6: Global velocity prediction normalized error (after the learning period) for the wake flow at a Reynolds number of 300 in the observation cases 2 (9 observation points, left) and 3 (1 observation point inside the recirculation zone, right) with, from top to bottom  $n = 2, 4$  and  $8$  modes for Red LUM (orange) and S-SOTA (light blue). The shaded colors (light orange and light blue) correspond to the respective estimated *posterior* standard deviations. The dashed black line at the bottom is the POD truncation error. The solid black line at the top is the error obtained by setting all temporal modes  $b_i$  to 0, i.e. keeping only the time averaged velocity  $\bar{v}$ .

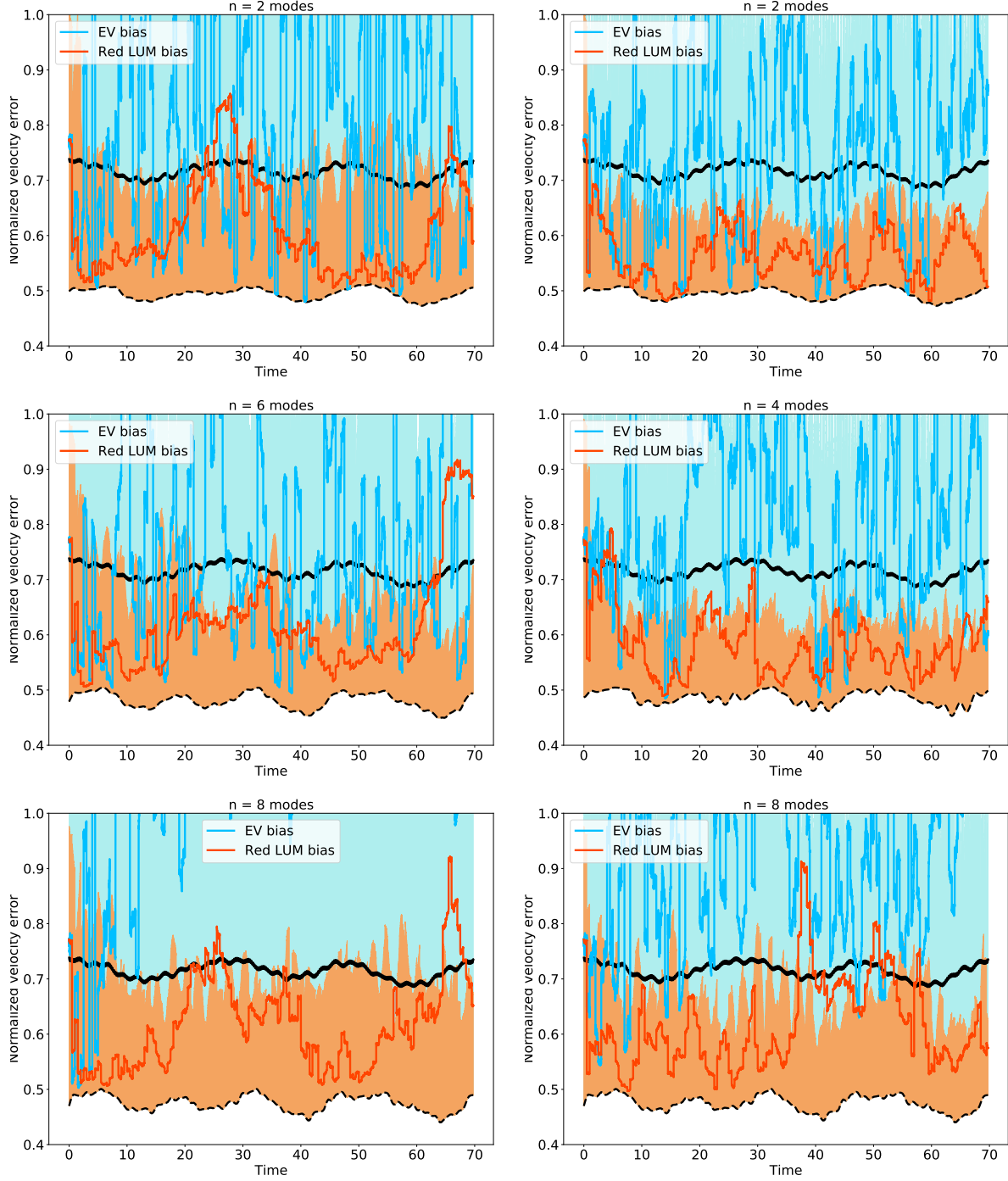


Figure 7: Global velocity prediction normalized error (after the learning period) for the wake flow at Reynolds 300 in the observation cases 4 (1 point outside the wake, left) and 5 (1 observation point far downstream, right) with, from top to bottom  $n = 2, 4$  and  $8$  modes for Red LUM (orange) and S-SOTA (light blue). The shaded colors (light orange and light blue) correspond to the respective estimated *posterior* standard deviations. The dashed black line at the bottom is the POD truncation error. The solid black line at the top is the error obtained by setting all temporal modes  $b_i$  to 0, i.e. keeping only the time averaged velocity  $\bar{v}$ .

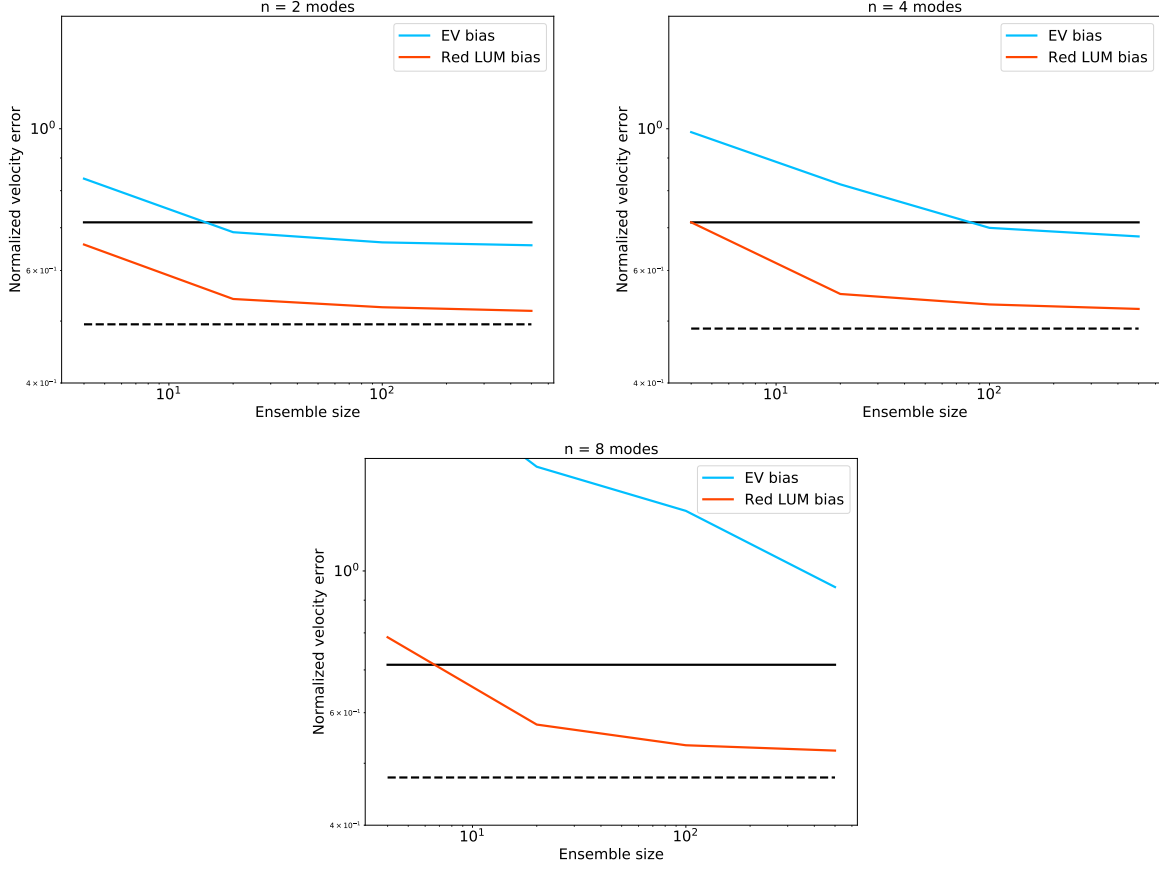


Figure 8: Loglog plots of the global velocity prediction normalized error (averaged on a time window of about 10 vortex shedding cycles, beginning about 2 vortex shedding cycles after the learning period) as a function of the ensemble size  $N_p$  for the wake flow at Reynolds 300 with, from left to right and from top to bottom  $n = 2, 4$  and 8 modes for Red LUM (orange) and S-SOTA (light blue). The dashed black line at the bottom is the POD truncation error. The solid black line at the top is the error obtained by setting all temporal modes  $b_i$  to 0, i.e. keeping only the time averaged velocity  $\bar{v}$ .

parallelized, most calculation run-times have been reduced to an order of minutes on a laptop, with the exception of the POD decomposition, and in particular the velocity temporal covariance calculation. This velocity temporal covariance calculation can be computationally demanding for large datasets (roughly  $10^4$  velocity snapshots of large dimensions in this instance).

The on-line data assimilation algorithm runs in a few minutes (about real-time) on a laptop computer, implemented without parallelization in Python. The time interval between two assimilation steps (0.5 s) is much larger than the time step used for the ROM time integration, and so the bulk of the computational expenditure is a product of the ROM time integration step. In the interval between consecutive assimilation steps, each realization is simulated independently. So, the CPU demand increases approximately linearly with the ensemble size, and moreover parallelization of the code could be achieved relatively easily. Led by the ROM time integration, the CPU consumption also scales with the number of coefficients in the ROM system, says  $O(n^3)$ . In practice, the ROM dimension  $n$ , can be selected to target a required accuracy. The eigenvalues of the POD decomposition express the amount of solution energy that can be encompassed for a specific value of  $n$ . A considerable further gains in run-time could be achieved through converting the online data assimilation code to C++, parallelizing the code, and through improvements to the time integration scheme. At term, the online algorithm should allow for measurements to be assimilated on the fly.

## 7. CONCLUSION

This work presents a novel reduced order model – referred to as Red LUM – to estimate and predict a flow velocity field in the entire domain in real time from sparse measurements. The proposed algorithm is based on a stochastic, low-dimensional system built from fundamental physics equations, reduced using simulated DNS data, and coupled with real-time measurements. The proposed approach has demonstrated much more accurate flow estimations, and convergence to an optimum solution with significantly smaller ensembles than current state-of-the-art reduced data assimilation methods, enabling a significant reduction in computational expenditure. Assimilating data from a single measurement point proved sufficient for our method to accurately predict the unsteady velocity field across the whole 3-dimensional domain. Combining particle filtering with models under location uncertainty, Red LUM has demonstrated an excellent performance in treating particle impoverishment. As such, Red LUM shows excellent promise for fluid dynamics applications requiring real-time fluid data assimilation.

Based on the initial promise of the proposed reduced order model, further work is planned to assimilate real, rather than simulated, particle image velocimetry data. Further enhancements are on going to adapt the model to the requirements of more turbulent flows thanks to the OpenFOAM-based Ithaca-FV library [52, 53]. Finally, it is envisaged to adapt Red LUM for parametric dependency and construction from noisy and possibly incomplete data sets.

## APPENDIX A. STOCHASTIC CALCULUS

Stochastic calculus is a field of Mathematics devoted to differential equations involving noises. Numerous tools exist in this framework, including theoretical moments computation, statistical estimations, or simulation of both ordinary differential equations [54] and partial differential equations [55, 56, 57, 58]. Several notations co-exist in this field, in particular Stratonovich or Itô conventions. This section introduces the Itô notation employed in this paper and some of its advantages. The discussion about relying whether on one notation or the other is recurrent in physics, but it is beyond the scope of this paper. For a more complete discussion on this subject, the reader can refer to sections 1.6 (pages 10-12) and 10.1 (pages 189-190) of [59]. For this short note, we highlight an important point: under appropriate assumptions, it is easy to switch from one notation to the other [54]. As a consequence, the most convenient form can be used to tackle a given issue.

In the Itô convention, time derivatives correspond to first-order forward-in-time differentials. Thus, the derivative,  $\partial_t w_k(\mathbf{x}, t)$  in equation (3) and  $\frac{db_i(t)}{dt}$  in equation (5) stand for  $\frac{w_k(\mathbf{x}, t+dt) - w_k(\mathbf{x}, t)}{dt}$  and  $\frac{b_i(t+dt) - b_i(t)}{dt}$  respectively, for an infinitesimal time step  $dt$ . The forward-in-time differential equations enable relatively simple computation. Alternative stochastic calculus conventions require much more complicate integration schemes. Furthermore, the Itô notation more explicitly identifies and separates the zero-mean noise terms  $-\mathbf{v}' \cdot \nabla w_k$  in equation (3) and  $\tilde{\alpha}_{pik}^R b_p(t) \dot{\beta}_k(t)$  in equation (5) – from the other terms induced by randomization of equations (e.g., diffusion, noise-induced drift). Consequently, Itô calculus greatly simplified the derivation of the evolution law of statistical moments (see for instance [27]).

## APPENDIX B. ESTIMATIONS OF THE REDUCED LOCATION UNCERTAINTY MODEL'S COEFFICIENTS

In this section, the various terms of the reduced location uncertainty model (5), initially proposed in a previous study by the current authors [29], are defined. In the following notation,  $(\boldsymbol{\zeta}, \boldsymbol{\xi}) \triangleq \int_{\Omega} \boldsymbol{\zeta} \cdot \boldsymbol{\xi}$  denotes the scalar product of the vectorial functions  $\boldsymbol{\zeta}$  and  $\boldsymbol{\xi}$ .  $\mathcal{P}$  denotes the non-local Leray operator  $\mathcal{P} = \mathbb{I}_d - \nabla \nabla^T \Delta^{-1}$ . The evaluation of this operator requires the resolution of a Poisson equation. It is used to simplify the fluid mechanics equations through the removal of the pressure term.

### Appendix B.1 Time down-sampling rate

Under the LU Navier-Stokes model hypothesis, the unresolved component of the velocity field  $\mathbf{v}'$  corresponds to a noise uncorrelated in time (i.e. white with respect to time). This assumption is consistent with the fact that the higher-order coefficients of the reduced-order solution often tend to have a shorter correlation time in fluid dynamics systems. However, in practice, this assumption has not proven to be entirely accurate and has presented recurring problems for the data-driven modeling of systems combining fast and slowly evolving components [60, 61, 62, 63].

Consequently, a time down-sampling scheme is used to force the noise terms to be as uncorrelated as possible.

Assuming that the spatially averaged covariance function has a Gaussian form with a standard deviation equal to the correlation time  $\tau$ , a simple expression allows the computation of the correlation time. For a given velocity correlation matrix  $C_{ij}^v = (\mathbf{v}_{\text{obs}}(\bullet, t_i), \mathbf{v}_{\text{obs}}(\bullet, t_j))$  (evaluated during the POD decomposition), the following  
 420 unresolved velocity correlation matrix can be evaluated:

$$C_{ij}^{v'} = (\mathbf{v}'_{\text{obs}}(\bullet, t_i), \mathbf{v}'_{\text{obs}}(\bullet, t_j)) = C_{ij}^v - \sum_{k=1}^n b_k^{\text{obs}}(t_i) b_k^{\text{obs}}(t_j), \quad 0 \leq i, j \leq N-1, \quad (\text{B.1})$$

along with its associated stationary covariance function

$$\text{Cov}_s(t_p) = \frac{1}{N-p} \sum_{q=0}^{N-1-p} C_{q, q+p}^{v'}, \quad 0 \leq p \leq N-1. \quad (\text{B.2})$$

The correlation time is estimated as follow:

$$\hat{\tau} = \sqrt{2 \frac{\overline{\text{Cov}_s^2}}{\left(\frac{\Delta \text{Cov}_s}{\Delta t}\right)^2}}, \quad (\text{B.3})$$

which is evaluated using a forward Euler temporal discretization of the stationary covariance:

$$\frac{\Delta \text{Cov}_s}{\Delta t}(t_p) \triangleq \frac{\text{Cov}_s(t_p + \Delta t) - \text{Cov}_s(t_p)}{\Delta t}, \quad 0 \leq p \leq N-1. \quad (\text{B.4})$$

Before computing the estimations presented in the following of this appendix, this estimated correlation time  $\hat{\tau}$   
 425 is used to down-sample the entire DNS velocity dataset and the observed coefficients of the reduced order solution, leaving us with a time step  $\Delta t \approx \hat{\tau}$ .

## Appendix B.2 Deterministic terms and noise-induced terms

The deterministic coefficients of Red LUM are summarized in table B.2.

## Appendix B.3 Noise correlations estimation

430 This section will discuss the estimation of the noise statistics.

For any function  $\xi$ , the linear functional  $K_{jq}$  can be defined as:

$$K_{jq}[\xi] \triangleq (\phi_j, -\mathcal{P}[(\xi \cdot \nabla) \phi_q] + \delta_{q0} \nu \Delta \xi), \quad 1 \leq j \leq n, \quad 0 \leq q \leq n. \quad (\text{B.5})$$

Using this notation, the noise's covariance can be estimated as follows:

$$\widehat{\Sigma_{pi, qj}^{\alpha}} = \frac{\Delta t}{\lambda_p^{\text{obs}}} K_{jq} \left[ \overline{b_p^{\text{obs}} \left( \frac{\Delta b_i^{\text{obs}}}{\Delta t} \right)'' v'_{\text{obs}}} \right], \quad 1 \leq i, j \leq n, \quad 0 \leq p, q \leq n, \quad (\text{B.6})$$



Physical meaning	Full-order term	ROM term
Molecular viscous dissipation	$\mathbf{L} = \frac{1}{Re} \Delta$	$l_{pi} = (\phi_i, \mathbf{L}(\phi_p))$
Usual advection	$\mathbf{C}(\mathbf{w}, \bullet) = -(\mathbf{w} \cdot \nabla)$	$c_{pqi} = (\phi_i, \mathcal{P}\mathbf{C}(\phi_p, \phi_q))$
Turbulent diffusion + Advecting velocity correction	$\mathbf{F}_{\text{dif}} = \nabla \cdot (\frac{1}{2} \mathbf{a} \nabla \bullet) + \frac{1}{2} (\nabla \cdot \mathbf{a}) \nabla$	$f_{pi} = (\phi_i, \mathcal{P}\mathbf{F}(\phi_p))$
Absolute diffusivity	$\mathbf{a}(\mathbf{x}) = \Delta t \overline{\mathbf{v}'_{\text{obs}} (\mathbf{v}'_{\text{obs}})^T}(\mathbf{x})$ $\mathbf{v}'_{\text{obs}} = \mathbf{v}_{\text{obs}} - \tilde{\Pi}_\phi[\mathbf{v}_{\text{obs}}]$	

Table B.2: Deterministic terms and noise-induced terms of the reduced location uncertainty model.

where  $b_0^{\text{obs}} = \lambda_0^{\text{obs}} = 1$  and for  $1 \leq i \leq n$ ,

$$b_i^{\text{obs}} = (\phi_i, \mathbf{v}_{\text{obs}}), \quad (\text{B.7})$$

$$\lambda_i^{\text{obs}} = \overline{(b_i^{\text{obs}})^2}, \quad (\text{B.8})$$

$$\left( \frac{\Delta b_i^{\text{obs}}}{\Delta t} \right)'' = \left( \frac{\Delta b_i^{\text{obs}}}{\Delta t} \right)' - \overline{\left( \frac{\Delta b_i^{\text{obs}}}{\Delta t} \right)'}, \quad (\text{B.9})$$

$$\left( \frac{\Delta b_i^{\text{obs}}}{\Delta t} \right)' = \left( \frac{\Delta b_i^{\text{obs}}}{\Delta t} \right) - \left( (\mathbf{b}^{\text{obs}})^T (\mathbf{l} + \mathbf{f})_{\bullet i} + (\mathbf{b}^{\text{obs}})^T \mathbf{c}_{\bullet \bullet i} \mathbf{b}^{\text{obs}} \right), \quad (\text{B.10})$$

$$\left( \frac{\Delta b_i^{\text{obs}}}{\Delta t} \right)(t_k) = \frac{b_i^{\text{obs}}(t_k + \Delta t) - b_i^{\text{obs}}(t_k)}{\Delta t}, \quad 0 \leq k \leq N-1. \quad (\text{B.11})$$

To ensure that noise covariance matrix conforms to the required symmetric, non-negative structure, the symmetric part of the estimated tensor (B.6) is retained and negative eigenvalues are set to zero. The consistency of this estimator is proven in a previous study by the current authors [29].

The dimension of the noise is reduced through a tensorial PCA of  $\widehat{\Sigma}^\alpha$ , retaining the initial  $n$  first eigenvectors.  $(\tilde{\alpha}_k^R)_{1 \leq k \leq n} \subset \mathbb{R}^{(n+1) \times n}$  are the matrix forms of the first  $n$  eigenvectors (weighted by the square roots of the corresponding eigenvalues). Since the noise is multiplicative and the temporal coefficients  $b_i$  have various amplitudes  $\sqrt{\lambda_i}$ , the covariance matrix  $\widehat{\Sigma}^\alpha$  is adequately re-normalized by the amplitudes  $\sqrt{\lambda_i}$  before applying the PCA.

## APPENDIX C. DESIGN OF THE OBSERVATION MODEL

An observation model aims at representing the link between measured values and an observed state. Here the state is  $\mathbf{b}$  and the observation model is represented in equation C.1

$$\mathbf{y} = \mathcal{H}[\mathbf{v}] + \epsilon_y = \underline{\mathcal{H}} \underline{\mathbf{v}} + \mathbf{L}_F \dot{\mathbf{W}}, \quad (\text{C.1})$$

where  $\mathbf{y}$  is the vector of the  $M_y$  PIV measurements to assimilate,  $\epsilon_y = \mathbf{L}_F \dot{\mathbf{W}}$  is the PIV measurement noise,  $\dot{\mathbf{W}}$  is a vector of  $M_{PIV}$  independent white noise (in discrete time),  $\mathbf{L}_F \mathbf{L}_F^T$  is the covariance matrix of the PIV measurement noise,  $\underline{\mathbf{v}}$  is the reshaped version of  $\mathbf{v}$  as a single vector of  $M \times d$  coefficients and  $\underline{\mathbf{H}}$  is a reshaped version of the linear operator  $\mathbf{H}$  as  $M_y \times (M \times d)$  matrix. For instance, if  $d = 3$ , we have  $\underline{\mathbf{v}}(t) =$   
445  $(v_1(\mathbf{x}_1, t), \dots, v_1(\mathbf{x}_M, t), v_2(\mathbf{x}_1, t), \dots, v_2(\mathbf{x}_M, t), v_3(\mathbf{x}_1, t), \dots, v_3(\mathbf{x}_M, t))^T$ . The two-scale representation  $\mathbf{v}(t) = \mathbf{w}(t) + \sigma \dot{\mathbf{B}}(t)$  can be extended to the reshaped version of  $\mathbf{v}$ :

$$\underline{\mathbf{v}} = \underline{\mathbf{w}} + \underline{\mathbf{v}}' = \underline{\Phi} \mathbf{b} + \underline{\sigma} \dot{\underline{\mathbf{B}}}, \quad (\text{C.2})$$

with  $\underline{\mathbf{w}}$  a reshaped vector of  $(M \times d)$  coefficients,  $\dot{\underline{\mathbf{B}}}$  is a vector of  $(M \times d)$  independent white noise (in discrete time),  $\underline{\sigma}$  a matrix of dimension  $(M \times d) \times (M \times d)$  and  $\underline{\Phi}$  the reshaped matrix  $\underline{\Phi} = (\underline{\phi}_0, \dots, \underline{\phi}_n)$  with dimension  $(M \times d) \times n$ . Finally, the observation model equation can be represented as the equation C.3

$$\mathbf{y} = \underline{\mathbf{H}} \underline{\Phi} \mathbf{b} + \underline{\mathbf{H}} \underline{\sigma} \dot{\underline{\mathbf{B}}} + \mathbf{L}_F \dot{\mathbf{W}}. \quad (\text{C.3})$$

450 The matrices  $\underline{\mathbf{H}}$  and  $\mathbf{L}_F$  are unknown and must be estimated. The  $\underline{\mathbf{H}}$  matrix is a spatial transformation from one space to another and mathematically it executes a spatial filtering and a cropping on the velocity field. We decompose the two operations through two operators:

$$\underline{\mathbf{H}} = \mathbf{H}_{\text{crop}} \mathbf{H}_{\text{PIV}}, \quad (\text{C.4})$$

where  $\mathbf{H}_{\text{crop}}$  is a rectangular matrix of 0 and 1 that defined the cropping of the PIV image and  $\mathbf{H}_{\text{PIV}}$  is a matrix representing the PIV measurement and postprocessing. The  $\mathbf{L}_F$  matrix is a quantification of uncertainty from  
455 the measures after post treatment. In order to estimate  $\mathbf{H}_{\text{PIV}}$  and  $\mathbf{L}_F$ , we have to compare the PIV data with a reference accurate measurement method : the hot wire.

Let us first focus on the matrix  $\mathbf{H}_{\text{PIV}}$ . Since we focus on 2D2C PIV, that matrix should include the slicing of the 3D data and the selection of the horizontal velocity components through two matrices of 0 and 1, that we name  
460  $\mathbf{H}_{\text{2D}}$  and  $\mathbf{H}_{\text{2C}}$  respectively:

$$\mathbf{H}_{\text{PIV}} = \mathbf{H}_{\text{2D}} \mathbf{H}_{\text{2C}} \mathbf{H}_{\text{Blur}}, \quad (\text{C.5})$$

with  $\mathbf{H}_{\text{Blur}}$  to encode the blurring induced by the PIV measurement. This blurring is approximated by an isotropic 3D spatial convolution on each velocity component:

$$v_k^{\text{PIV}} = h^S * v_k + \epsilon_k^{\text{PIV}}, \quad (\text{C.6})$$

where  $h^S$  is the spatial filter and  $\epsilon^{\text{PIV}}$  the measurement noise. We consider an isotropic Gaussian filter.

$$h^S(x, y, z) = h^s(x) h^s(y) h^s(z), \quad (\text{C.7})$$

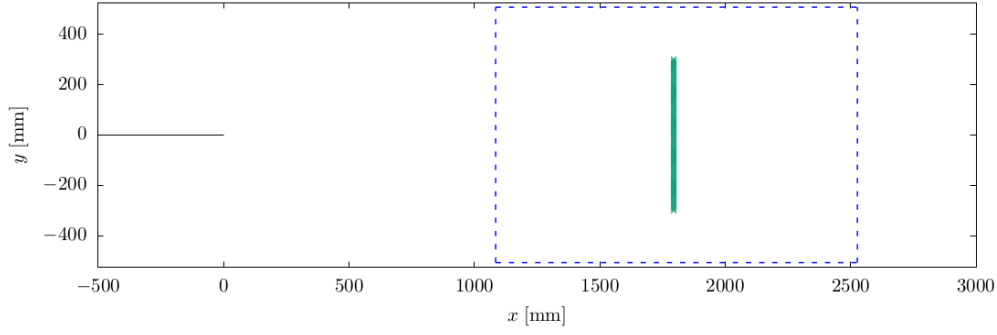


Figure C.9: Measure configuration plan of PIV and hot-wire.[64]

with

$$h^s(x) = A_h \exp \left( -\frac{1}{2} \left( \frac{x}{\sigma_s^h} \right)^2 \right). \quad (\text{C.8})$$

In Fourier space, the modulus of the filter is also a Gaussian function and can then be estimated as:

$$|\hat{h}^S| \approx \frac{|\hat{v}_k^{\text{PIV}}|}{|\hat{v}_k|}. \quad (\text{C.9})$$

Unfortunately, the hot wire data do not have a significant spatial extension. Accordingly, for our experimental data comparison we must focus on the temporal signatures. We will estimate a blurring temporal filter  $h^t$  instead of the blurring spatial filter  $h^s$ . And subsequently, we will transform the temporal filter into a spatial filter using a Taylor assumption. Let us now estimate the blurring temporal filter  $h^t$  from the experimental, keeping a Gaussian filter form.

Since the large-scale structures should be well reconstructed by the PIV, we can expect that the first Fourier mode is  $\hat{h}^t(0) = 1$ . This sets the amplitude of the filter. In Fourier space, the modulus of the filter is Gaussian and its logarithm writes:

$$\log \left( \frac{|\hat{v}_k^{\text{PIV}}(f)|}{|\hat{v}_k(f)|} \right) \approx \log \left( |\hat{h}^t(f)| \right) = -\alpha f^2. \quad (\text{C.10})$$

This gives a simple way to fit the last filter parameter: the standard deviation  $\sigma_t^h = \sqrt{\frac{\alpha}{2\pi^2}}$ . It is estimated with a simple linear regression in loglog plot.

The data used to estimate the matrices were presented in the work of [64]. The measure configuration is detailed by the Figure C.9, where the blue dotted line bounds the PIV area estimation and the green line represents the hot wire converters position. The sampling frequency used in PIV is 500 Hertz and 6000 Hertz in hot wire. The hot wire data spectrum and the PIV data spectrum are compared in figure C.10. The Gaussian fitting of  $\hat{h}^t$  is illustrated in figure C.11.

After that, the Taylor hypothesis is used again to transform the time filter  $h^t$  into a space filter  $h^s$  by the axis rescaling  $x = Ut$  where  $U$  is the mean flow velocity. Therefore, the space filter  $h^s$  is a Gaussian function with a

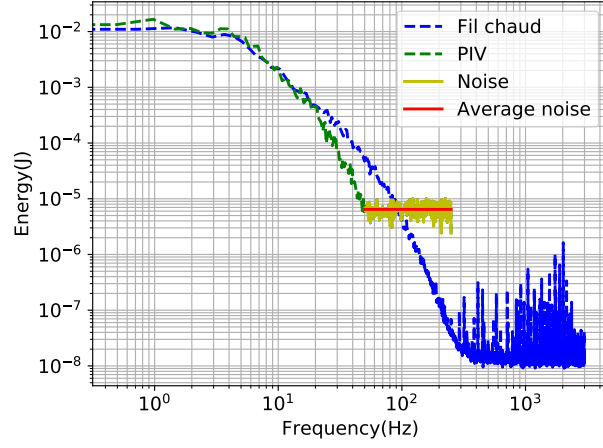


Figure C.10: Spectral analysis of PIV measurements versus hot-wire measurements.

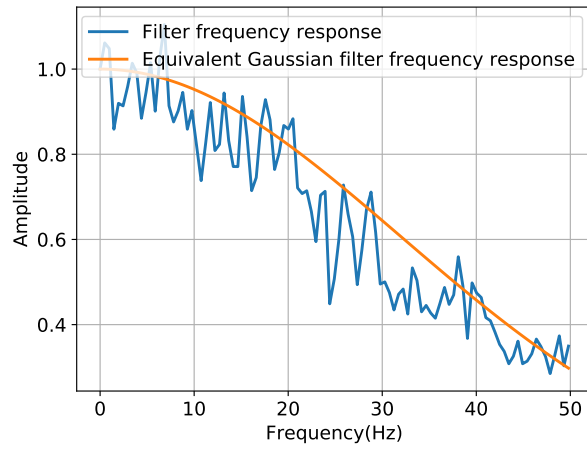


Figure C.11: Gaussian fitting of  $|\hat{h}^t|$ .

standard deviation  $\sigma_x^h = U\sigma_t^h$ .

Eventually, the filter  $h^s$  is used to define the PIV blurring matrix  $\mathbf{H}_{\text{Blur}}$ . According to equation (C.7), we can write the matrix  $\mathbf{H}_{\text{PIV}}$  as a reshaped version of a tensor product of 3 matrices representing the spatial smoothing along  $x$ ,  $y$  and  $z$ :

$$\mathbf{H}_{\text{Blur}} = \underline{\mathbf{H}_{\text{PIV}}^x} \otimes \underline{\mathbf{H}_{\text{PIV}}^y} \otimes \underline{\mathbf{H}_{\text{PIV}}^z}. \quad (\text{C.11})$$

Each of this matrix is defined in the same way though the filter  $h^s$ . In order to restrict the number of non-zero coefficient of those huge matrices, we crop the filter after one standard deviation  $h(x) = h^s(x)\mathbb{1}_{\{|x| \leq \sigma_x\}}$ . and

$$\mathbf{H}_{\text{PIV}}^x = \begin{bmatrix} h(-n_h\Delta x) & \dots & h(0) & \dots & h(n_h\Delta x) & \dots & \dots & \dots & 0 \\ 0 & \dots & \dots & \dots & \dots & \dots & \dots & \dots & 0 \\ 0 & \dots & h(-n_h\Delta x) & \dots & h(0) & \dots & h(n_h\Delta x) & \dots & 0 \\ 0 & \dots & \dots & \dots & \dots & \dots & \dots & \dots & 0 \\ 0 & \dots & \dots & \dots & h(-n_h\Delta x) & \dots & h(0) & \dots & h(n_h\Delta x) \end{bmatrix}, \quad (\text{C.12})$$

where  $n_h$  is the floor function of  $\sigma_x/\Delta x$ .

This finally leads to the observation matrix

$$\underline{\mathbf{H}} = \mathbf{H}_{\text{crop}} \mathbf{H}_{\text{PIV}} = \mathbf{H}_{\text{crop}} \mathbf{H}_{2\text{D}} \mathbf{H}_{2\text{C}} \underline{\mathbf{H}_{\text{PIV}}^x} \otimes \underline{\mathbf{H}_{\text{PIV}}^y} \otimes \underline{\mathbf{H}_{\text{PIV}}^z}. \quad (\text{C.13})$$

On top of the observation matrix, we need to estimate the matrix  $\mathbf{L}_F$ . To do so, we use the noisy part of the PIV spectrum (see figure C.10). To simplify we assume that the noise  $\epsilon_y$  is white in time and in space. We estimate the noise variance  $\sigma_\epsilon^2$  from the PIV spectrum. With the cropping of this PIV image, this leads to:

$$\mathbf{L}_{PIV} = \mathbf{H}_{\text{crop}} (\sigma_\epsilon \mathbb{I}_d) = \sigma_\epsilon \mathbf{H}_{\text{crop}}. \quad (\text{C.14})$$

## APPENDIX D. LOG-LIKELIHOOD EXPRESSION

The log-likelihood is necessary in order to assimilate the PIV measurements. Since the observation model (C.3) is linear with an additive Gaussian noise  $\underline{\epsilon}_y^R = \underline{\mathbf{H}} \underline{\boldsymbol{\sigma}} \underline{\dot{\mathbf{B}}} + \mathbf{L}_F \dot{\mathbf{W}}$ , the log-likelihood is a quadratic function of the state  $\mathbf{b}$ :

$$p(\mathbf{y}(t)|\mathbf{b}(t)) \propto \exp\left(-\frac{1}{2}\|\mathbf{y}(t) - \underline{\mathbf{H}} \underline{\boldsymbol{\Phi}} \mathbf{b}(t)\|_{\boldsymbol{\Sigma}^{-1}}^2\right), \quad (\text{D.1})$$

$$\propto \exp\left(-\frac{1}{2}\left(\mathbf{y}^T(t)\boldsymbol{\Sigma}^{-1}\mathbf{y}(t) + \mathbf{y}^T(t)\mathbf{B}\mathbf{b}(t) + \mathbf{b}^T(t)\mathbf{A}\mathbf{b}(t)\right)\right), \quad (\text{D.2})$$

with  $\boldsymbol{\Sigma}$  the covariance matrix of the whole additive noise  $\underline{\epsilon}_y^R = \underline{\mathbf{H}} \underline{\boldsymbol{\sigma}} \underline{\dot{\mathbf{B}}} + \mathbf{L}_F \dot{\mathbf{W}}$

$$\boldsymbol{\Sigma} = (\underline{\mathbf{H}} \underline{\boldsymbol{\sigma}}) (\underline{\mathbf{H}} \underline{\boldsymbol{\sigma}})^T + \mathbf{L}_F \mathbf{L}_F^T, \quad (\text{D.3})$$

490 and

$$\mathbf{A} = -\frac{1}{2}(\underline{\mathbf{H}} \underline{\Phi})^T \mathbf{B}, \quad (\text{D.4})$$

$$\mathbf{B} = \Sigma^{-1}(\underline{\mathbf{H}} \underline{\Phi}). \quad (\text{D.5})$$

In practice,  $\underline{\sigma}$  and  $\underline{\sigma} \underline{\sigma}^T$  are huge and cannot be even memorize. Depending on the PIV cropping,  $\Sigma$  can be huge as well which prevents the computation of its inverse. Therefore, for the computation of  $\Sigma$  and its inverse only, we estimate  $(\underline{\mathbf{H}} \underline{\sigma})(\underline{\mathbf{H}} \underline{\sigma})^T$  with the following time average:

$$(\underline{\mathbf{H}} \underline{\sigma})(\underline{\mathbf{H}} \underline{\sigma})^T = \overline{(\mathbf{H}[v']) (\mathbf{H}[v'])^T} = \mathbf{H}_{\text{crop}} \overline{(\mathbf{H}_{\text{PIV}} \underline{v}') (\mathbf{H}_{\text{PIV}} \underline{v}')^T} \mathbf{H}_{\text{crop}}^T. \quad (\text{D.6})$$

495 Further simplification is achieved through the negations of the spatial correlation in  $\mathbf{H}[v']$ . Subsequently,  $(\underline{\mathbf{H}} \underline{\sigma})(\underline{\mathbf{H}} \underline{\sigma})^T$  and then  $\Sigma$  become quasi-diagonal, which enables the computation of  $\Sigma^{-1}$ .

Besides, the factor  $\exp(-\frac{1}{2}\mathbf{y}(t)^T \Sigma^{-1} \mathbf{y}(t))$  in equation (D.2) is equivalent for all realizations  $\mathbf{b}^{(j)}$ . So, it will disappear in the state distribution normalization step and so does not required computation:

$$p(\mathbf{y}(t)|\mathbf{b}(t)) \propto \exp\left(\mathbf{y}^T(t) \mathbf{B} \mathbf{b}(t) + \mathbf{b}^T(t) \mathbf{A} \mathbf{b}(t)\right). \quad (\text{D.7})$$

## SAMPLE CREDIT AUTHOR STATEMENT

**Valentin Resseguier:** Conceptualization, Methodology, Software, Visualization, Writing. **Matheus Ladvig:** Software, Investigation, Formal analysis. **Dominique Heitz:** Data Curation, Methodology.

## ACKNOWLEDGMENTS

500 We warmly thank Michael Jonhson for a careful proofreading of this paper and its many corrections, Pranav Chandramouli for the generation of the three-dimensional wake flow direct numerical simulation (DNS) data and Romain Schuster for providing the experimental data. We also thank Darryl D. Holm, Dan Crisan, Igor Shevchenko and above all Wei Pan for the insightful discussions.

## FUNDING

505 This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors. Two of the authors have been employed by a private company (Scalian) when performing this work.

## REFERENCES

- [1] D. Schlipf, Lidar-assisted control concepts for wind turbines, Ph.D. thesis (2016).

- [2] A. Soulier, C. Braud, D. Voisin, B. Podvin, Low-reynolds-number investigations on the ability of the strip of e-telltale sensor to detect the flow features over wind turbine blade section: flow stall and reattachment dynamics, *Wind Energy Science* 6 (2) (2021) 409–426.
- 510 [3] E. Livne, Aircraft active flutter suppression: State of the art and technology maturation needs, *Journal of Aircraft* 55 (1) (2018) 410–452.
- [4] N. B. Erichson, L. Mathelin, Z. Yao, S. L. Brunton, M. W. Mahoney, J. N. Kutz, Shallow neural networks for fluid flow reconstruction with limited sensors, *Proceedings of the Royal Society A* 476 (2238) (2020) 20200097.
- 515 [5] J. P. Thomas, E. H. Dowell, K. C. Hall, Three-dimensional transonic aeroelasticity using proper orthogonal decomposition-based reduced-order models, *Journal of Aircraft* 40 (3) (2003) 544–551.
- [6] C. Braud, D. Heitz, G. Arroyo, L. Perret, J. Deleville, J. Bonnet, Low-dimensional analysis, using POD, for two mixing layer-wake interactions, *International Journal of Heat and Fluid Flow* 3 (3) (2004) 351–363.
- [7] L. Fick, Y. Maday, A. T. Patera, T. Taddei, A stabilized pod model for turbulent flows over a range of reynolds numbers: Optimal parameter sampling and constrained projection, *Journal of Computational Physics* 371 (2018) 214–243.
- 520 [8] A. Majda, Statistical energy conservation principle for inhomogeneous turbulent dynamical systems, *Proceedings of the National Academy of Sciences* 112 (29) (2015) 8937–8941.
- [9] T. Sapsis, Attractor local dimensionality, nonlinear energy transfers and finite-time instabilities in unstable dynamical systems with applications to two-dimensional fluid flows, *Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences* 469 (2153). [doi:10.1098/rspa.2012.0550](https://doi.org/10.1098/rspa.2012.0550).
- 525 [10] T. Sapsis, A. Majda, Blending modified Gaussian closure and non-Gaussian reduced subspace methods for turbulent dynamical systems, *Journal of Nonlinear Science* 23 (6) (2013) 1039–1071.
- [11] N. Aubry, P. Holmes, J. Lumley, E. Stone, The dynamics of coherent structures in the wall region of a turbulent boundary layer, *J. Fluid Mech.* 192 (1988) 115–173.
- 530 [12] W. Cazemier, R. Verstappen, A. Veldman, Proper orthogonal decomposition and low-dimensional models for driven cavity flows, *Phys. Fluids* 10 (7) (1998) 1685–1699.
- [13] Z. Wang, I. Akhtar, J. Borggaard, T. Iliescu, Proper orthogonal decomposition closure models for turbulent flows: a numerical comparison, *Computer Methods in Applied Mechanics and Engineering* 237 (2012) 10–26.
- 535 [14] M. Buffoni, S. Camarri, A. Iollo, M. V. Salvetti, Low-dimensional modelling of a confined three-dimensional wake flow, *Journal of Fluid Mechanics* 569 (2006) 141–150.

- [15] L. Cordier, B. R. Noack, G. Tissot, G. Lehnasch, J. Delville, M. Balajewicz, G. Daviller, R. K. Niven, Identification strategies for model-based control, *Experiments in fluids* 54 (8) (2013) 1580.
- [16] X. Xie, M. Mohebbujjaman, L. G. Rebholz, T. Iliescu, Data-driven filtered reduced order modeling of fluid flows, *SIAM Journal on Scientific Computing* 40 (3) (2018) B834–B857.
- 540 [17] I. M. Navon, Data assimilation for numerical weather prediction: a review, in: *Data assimilation for atmospheric, oceanic and hydrologic applications*, Springer, 2009, pp. 21–65.
- [18] M. Couplet, C. Basdevant, P. Sagaut, Calibrated reduced-order pod-galerkin system for fluid flow modelling, *Journal of Computational Physics* 207 (1) (2005) 192–220.
- [19] J. D’adamo, N. Papadakis, E. Memin, G. Artana, Variational assimilation of pod low-order dynamical systems, *Journal of Turbulence* 8 (9) (2007) 1–22.
- 545 [20] G. Artana, A. Cammilleri, J. Carlier, E. Mémin, Strong and weak constraint variational assimilations for reduced order fluid flow modeling, *J. Comp. Phys* 231 (8) (2012) 3264–3288.
- [21] R. Kikuchi, T. Misaka, S. Obayashi, International journal of computational fluid dynamics real-time prediction of unsteady flow based on pod reduced-order model and particle filter, *International Journal of Computational*
- 550 *Fluid Dynamics* 30 (4) (2016) 285–306.
- [22] E. Mémin, Fluid flow dynamics under location uncertainty, *Geophysical & Astrophysical Fluid Dynamics* 108 (2) (2014) 119–146.
- [23] V. Resseguier, E. Mémin, B. Chapron, Geophysical flows under location uncertainty, part I random transport and general models, *Geophysical & Astrophysical Fluid Dynamics* 111 (3) (2017) 149–176.
- 555 [24] R. Mikulevicius, B. Rozovskii, Stochastic Navier–Stokes equations for turbulent flows, *SIAM Journal on Mathematical Analysis* 35 (5) (2004) 1250–1310.
- [25] V. Resseguier, E. Mémin, B. Chapron, Geophysical flows under location uncertainty, part II quasi-geostrophy and efficient ensemble spreading, *Geophysical & Astrophysical Fluid Dynamics* 111 (3) (2017) 177–208.
- [26] V. Resseguier, W. Pan, B. Fox-Kemper, Data-driven versus self-similar parameterizations for stochastic advection by lie transport and location uncertainty, *Nonlinear Processes in Geophysics* 27 (2) (2020) 209–234.
- 560 [27] V. Resseguier, L. Li, G. Jouan, P. Dérian, E. Mémin, C. Bertrand, New trends in ensemble forecast strategy: uncertainty quantification for coarse-grid computational fluid dynamics, *Archives of Computational Methods in Engineering* (2020) 1–82.



- [28] B. Chapron, P. Dérian, E. Mémin, V. Resseguier, Large-scale flows under location uncertainty: a consistent stochastic framework, *Quarterly Journal of the Royal Meteorological Society* 144 (710) (2018) 251–260.
- [29] V. Resseguier, A. M. Picard, E. Memin, B. Chapron, Quantifying truncation-related uncertainties in unsteady fluid dynamics reduced order models, *SIAM/ASA Journal on Uncertainty Quantification* 9 (3) (2021) 1152–1183.
- [30] W. Bauer, P. Chandramouli, L. Li, E. Mémin, Stochastic representation of mesoscale eddy effects in coarse-resolution barotropic models, *Ocean Modelling* 151 (2020) 101646.
- [31] R. Brecht, L. Li, W. Bauer, E. Mémin, Rotating shallow water flow under location uncertainty with a structure-preserving discretization, *Journal of Advances in Modeling Earth Systems* 13 (12) (2021) e2021MS002492.
- [32] V. Resseguier, E. Mémin, D. Heitz, B. Chapron, Stochastic modelling and diffusion modes for proper orthogonal decomposition models and small-scale flow analysis, *Journal of Fluid Mechanics* 826 (2017) 888–917.
- [33] A. Doucet, A. Johansen, A tutorial on particle filtering and smoothing: Fifteen years later, *Handbook of Nonlinear Filtering* 12 (2009) 656–704.
- [34] A. Doucet, N. De Freitas, N. Gordon, *Sequential Monte Carlo methods in practice*, Springer, 2001.
- [35] N. Kantas, A. Beskos, A. Jasra, Sequential monte carlo methods for high-dimensional inverse problems: A case study for the navier–stokes equations., *SIAM/ASA Journal on Uncertainty Quantification* 2.1 (2014) 464–489.
- [36] P. Rebeschini, R. Van Handel, et al., Can local particle filters beat the curse of dimensionality?, *The Annals of Applied Probability* 25 (5) (2015) 2809–2866.
- [37] A. Beskos, D. Crisan, A. Jasra, K. Kamatani, Y. Zhou, A stable particle filter for a class of high-dimensional state-space models, *Advances in Applied Probability* 49 (1) (2017) 24–48.
- [38] C. Cotter, D. Crisan, D. D. Holm, W. Pan, I. Shevchenko, A particle filter for stochastic advection by lie transport: A case study for the damped and forced incompressible two-dimensional euler equation, *SIAM/ASA Journal on Uncertainty Quantification* 8 (4) (2020) 1446–1492.
- [39] A. Farchi, M. Bocquet, Comparison of local particle filters and new implementations., *Nonlinear Processes in Geophysics* 25 (4).
- [40] F.-X. Le Dimet, O. Talagrand, Variational algorithms for analysis and assimilation of meteorological observations: theoretical aspects, *Tellus A* 38 (2) (1986) 97–110.
- [41] J. S. Liu, *Monte Carlo strategies in scientific computing*, Springer Science & Business Media, 2008.

- [42] A. Doucet, X. Wang, Monte carlo methods for signal processing: a review in the statistical signal processing context, *IEEE Signal Processing Magazine* 22 (6) (2005) 152–170.
- [43] S. Laizet, E. Lamballais, High-order compact schemes for incompressible flows: a simple and efficient method with the quasi-spectral accuracy, *J. Comp. Phys.* 228 (15) (2009) 5989–6015.
- [44] J. Berner, S.-Y. Ha, J. Hacker, A. Fournier, C. Snyder, Model uncertainty in a mesoscale ensemble prediction system: Stochastic versus multiphysics representations, *Monthly Weather Review* 139 (6) (2011) 1972–1995.
- [45] C. Franzke, T. O’Kane, J. Berner, P. Williams, V. Lucarini, *Stochastic climate theory and modeling*, Wiley Interdisciplinary Reviews: Climate Change 6 (1) (2015) 63–78.
- [46] G. Gottwald, J. Harlim, The role of additive and multiplicative noise in filtering complex dynamical systems, *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Science* 469 (2155) (2013) 20130096.
- [47] L. Mitchell, G. Gottwald, Data assimilation in slow-fast systems using homogenized climate models, *Journal of the atmospheric sciences* 69 (4) (2012) 1359–1377.
- [48] C. Penland, L. Matrosova, A balance condition for stochastic numerical models with application to the El Nino-southern oscillation, *Journal of climate* 7 (9) (1994) 1352–1372.
- [49] C. Penland, P. Sardeshmukh, The optimal growth of tropical sea surface temperature anomalies, *Journal of climate* 8 (8) (1995) 1999–2024.
- [50] A. Trevisan, F. Uboldi, Assimilation of standard and targeted observations within the unstable subspace of the observation-analysis-forecast cycle system, *Journal of the atmospheric sciences* 61 (1) (2004) 103–113.
- [51] G.-H. Ng, D. McLaughlin, D. Entekhabi, A. Ahanin, The role of model dynamics in ensemble Kalman filter performance for chaotic systems, *Tellus A* 63 (5) (2011) 958–977.
- [52] G. Stabile, S. Hijazi, A. Mola, S. Lorenzi, G. Rozza, POD-Galerkin reduced order methods for CFD using finite volume discretisation: vortex shedding around a circular cylinder, *Communications in Applied and Industrial Mathematics* 8 (1) (2017) 210–236.
- [53] G. Stabile, G. Rozza, Finite volume POD-Galerkin stabilised reduced order methods for the parametrised incompressible navier–stokes equations, *Computers & Fluids* 173 (2018) 273–284.
- [54] B. Oksendal, *Stochastic differential equations*, Springer-Verlag, 1998.
- [55] H. Kunita, *Stochastic flows and stochastic differential equations*, Vol. 24, Cambridge university press, 1997.

- [56] G. Da Prato, J. Zabczyk, Stochastic Equations in Infinite Dimensions, Encyclopedia of Mathematics and its Applications, Cambridge University Press, 1992.
- [57] C. Prévôt, M. Röckner, A concise course on stochastic partial differential equations, Vol. 1905, Springer, 2007.
- [58] M. Choi, T. Sapsis, G. Karniadakis, On the equivalence of dynamically orthogonal and bi-orthogonal methods: Theory and numerical simulations, Journal of Computational Physics 270 (2014) 1–20.
- [59] V. Resseguier, Mixing and fluid dynamics under location uncertainty, Ph.D. thesis, Rennes 1 university (2017).
- [60] R. Azencott, A. Beri, A. Jain, I. Timofeyev, Sub-sampling and parametric estimation for multiscale dynamics, Communications in Mathematical Sciences 11 (4) (2013) 939–970.
- [61] R. Azencott, A. Beri, I. Timofeyev, Adaptive sub-sampling for parametric estimation of Gaussian diffusions, Journal of Statistical Physics 139 (6) (2010) 1066–1089.
- [62] A. Papavasiliou, G. Pavliotis, A. Stuart, Maximum likelihood drift estimation for multiscale diffusions, Stochastic Processes and their Applications 119 (10) (2009) 3173–3210.
- [63] G. Pavliotis, A. Stuart, Parameter estimation for multiscale diffusions, Journal of Statistical Physics 127 (4) (2007) 741–781.
- [64] R. Schüster, Développement d’une méthode de mesure basée image pour caractériser en grande taille les flux d’air intérieurs, Ph.D. thesis, Université de Rennes I (2019).