



Adaptive ensemble optimal interpolation for efficient data assimilation in the red sea

Item Type	Article
Authors	Toye, Habib;Zhan, Peng;Sana, Furrukh;Sanikommu, Siva Reddy;Raboudi, Naila Mohammed Fathi;Hoteit, Ibrahim
Citation	Toye, H., Zhan, P., Sana, F., Sanikommu, S., Raboudi, N., & Hoteit, I. (2021). Adaptive ensemble optimal interpolation for efficient data assimilation in the red sea. Journal of Computational Science, 51, 101317. doi:10.1016/j.jocs.2021.101317
Eprint version	Post-print
DOI	10.1016/j.jocs.2021.101317
Publisher	Elsevier BV
Journal	Journal of Computational Science
Rights	NOTICE: this is the author's version of a work that was accepted for publication in Journal of Computational Science. Changes resulting from the publishing process, such as peer review, editing, corrections, structural formatting, and other quality control mechanisms may not be reflected in this document. Changes may have been made to this work since it was submitted for publication. A definitive version was subsequently published in Journal of Computational Science, [51, , (2021-02-06)] DOI: 10.1016/j.jocs.2021.101317 . © 2021. This manuscript version is made available under the CC-BY-NC-ND 4.0 license http://creativecommons.org/licenses/by-nc-nd/4.0/
Download date	2024-03-13 10:05:30
Link to Item	http://hdl.handle.net/10754/667512

Adaptive Ensemble Optimal Interpolation for Efficient Data Assimilation in the Red Sea

Habib Toye^a, Peng Zhan^a, Furruxh Sana^{a,b}, Sivareddy Sanikommu^a, Naila Raboudi^a, Ibrahim Hoteit^{a,*}

^a *King Abdullah University of Science and Technology, Saudi Arabia*

^b *Harvard Medical School, Massachusetts General Hospital, USA*

Abstract

Ensemble optimal interpolation (EnOI) is a variant of the ensemble Kalman filter (EnKF) that operates with a static ensemble to drastically reduce its computational cost. The idea is to use a pre-selected ensemble to parameterize the background covariance matrix, which avoids the costly integration of the ensemble members with the dynamical model during the forecast step of the filtering process. To better represent the pronounced time-varying circulation of the Red Sea, we propose a new adaptive EnOI approach in which the ensemble members are adaptively selected at every assimilation cycle from a large dictionary of ocean states describing the Red Sea variability. We implement and test different schemes to select the ensemble members (i) based on the similarity to the forecast state according to some criteria, or (ii) in term of best representation of the forecast in an ensemble subspace using an Orthogonal Matching Pursuit (OMP) algorithm. The relevance of the schemes is first demonstrated with the Lorenz 63 and Lorenz 96 models. Then results of numerical experiments assimilating real remote sensing data into a high resolution MIT general circulation model (MITgcm) of the Red Sea using the Data Assimilation Research Testbed (DART) system are presented and discussed.

Keywords: Red Sea, Data Assimilation, Ensemble Kalman Filter, Ensemble Optimal Interpolation, Orthogonal Matching Pursuit

*Corresponding author

Email address: ibrahim.hoteit@kaust.edu.sa (Ibrahim Hoteit)

1. Introduction

The Red Sea lies between Africa and the Arabian Peninsula, connecting the Mediterranean Sea to the Indian Ocean. It serves chief shipping trade routes between Europe and Asia with the opening of the Suez Canal and substantially contributes to the social and economic developments in the surrounding countries. The Red Sea's complex terrains and landforms have accommodated a unique ecological system rich of biodiversity [4]. The Red Sea hosts a dozens of islands scattered along its coast, particularly in the Southern basin, while flat coastal plains gradually slope to the central trough, deepening up to 2100 m. The water exchange between the Red Sea and the Mediterranean Sea is quite limited, but the outflow/intrusion to/from the Gulf of Aden is significant for the Red Sea circulation and ecosystem despite the topographic restrictions of the narrow strait of Bab-Al-Mandeb [9, 65, 66].

An increasing number of observations have become available in recent years in the Red Sea, collected by different platforms, including satellites, buoys, gliders, cruises, etc. Despite this, the spatial and temporal data coverage in the Red Sea is still limited, making it difficult to study the hydrodynamics based solely on observations. Numerical simulations of the Red Sea circulation using eddy-resolving Oceanic General Circulation Models (OGCM) have therefore become popular to conduct various studies of the Red Sea circulation, including the general and overturning circulations [65, 66], deep water formations [46, 64], eddies variability and dynamics [67, 68, 69, 70], connectivity [40, 48], internal waves [17, 18], etc. Such models are also the step-stone to develop forecasting capabilities.

Numerical models are imperfect and inevitable source of errors may affect their outputs [11]. Currently, the best approach to obtain accurate model simulations is to condition the model outputs to available observations through a data assimilation process. Data assimilation methods seek for the best possible estimates of the state of a dynamical system given available data [21]. These are now widely used in atmospheric and oceanic applications for operational services, and also for parameter estimation, optimal observations design, etc [20, 52, 67].

Data assimilation in marginal seas is more challenging than in the open ocean, due to the rich nonlinear dynamics of coastal regions, the lack of coordinated observational networks, albeit for targeted efforts to monitor some specific local features of interest, and inevitably modeling uncertainties resulting from the complex topography, and coarse atmospheric forcing fields and ocean boundaries [11]. Until very recently, data assimilation was

not applied in the Red Sea, except for one early work [5] that implemented a coarse model with simplified atmospheric forcing and assimilated sea surface temperature (SST) and XBT/CTD data using a simple nudging technique. With the ever increasing interests to develop modeling and forecasting capabilities for the Red Sea, driven by the new mega Saudi governmental projects along the shores of the Red Sea and the desire of the national oil company ARAMCO to develop advanced operational capabilities to support their operations, significant efforts are now being undertaken to develop such a system using state-of-the-art modeling and data assimilation techniques. These were initiated by the Red Sea Modeling and Prediction Group at King Abdullah University of Science and Technology (KAUST). In 2017, [60] presented the implementation of the first Red Sea ensemble assimilation and forecasting system and assessed the performances of different assimilation schemes, namely a deterministic Ensemble Adjustment Kalman Filter (EAKF), and an Ensemble Optimal Interpolation (EnOI) with pre-selected static and seasonally varying ensembles. Composed of the MIT general circulation model (MITgcm) for ocean forecasting and the Data Assimilation Research Testbed (DART) for ensemble assimilation, the system is forced with realistic atmospheric fields and assimilates remotely sensed sea surface height (SSH) and SST observations. The system is further equipped to assimilate most, if not all, available ocean measurements [20, 60].

Ensemble assimilation methods have been proven very efficient in many ocean applications and regions (e.g. [6, 20, 43, 60, 63]). The performance of these methods greatly depends on the representativeness of their ensembles, which should be large enough to describe the directions of errors growth of the system and mitigate the effects of sampling errors [25, 32, 56]. Using large ensembles in an ensemble Kalman filter (EnKF), such as EAKF, means more numerical model integrations and therefore increased computational cost [60]. EnOI integrates only the filter estimate (i.e. analysis) to compute the forecast and updates the latter with the incoming observations based on the sample covariance of a pre-selected ensemble, as a way to reduce the number of model runs. A stationary ensemble may however not properly capture the striking seasonal variability of the Red Sea dynamics [65, 66]. A Seasonal EnOI, which uses seasonally varying ensembles [63], was successfully implemented in the Red Sea [23, 60]. It was however not very efficient at describing the prevailing eddy and mesoscale activities in the basin [69].

The use of pre-selected time-varying ensembles that represent the different seasons of the studied basin has already been proposed in EnOI [62, 63]. Here we propose to push this idea further by adaptively and automatically selecting, at every filter analysis step, a new ensemble from an available

“dictionary” of representative ocean states (e.g. long reanalysis). As in an EnKF, this will enable updating the ensemble of the EnOI scheme, in order to describe the state uncertainties at the time of the analysis step, while avoiding the costly numerical integration of its members. The selection of the ensemble members will be based on the best estimate of the state at the time of the analysis, i.e. the forecast state. This new Adaptive EnOI (AEnOI) scheme will therefore only integrate the model once to forecast the state, and then select an ensemble from the dictionary that represents its current uncertainties according to a certain criteria, based on which the Kalman analysis step will be applied.

Similar ideas have been recently proposed, relying on some kind of dictionary to describe the uncertainties or statistics of the estimate of interest. [58], for instance, suggested a fully data-driven ensemble data assimilation framework that selects the “best” ensemble members from a given “catalog” of possible successive states of the system based on a “analog” or “nearest-neighbor” approach. The nearest-neighbor approach is adopted from the machine learning community and is basically designed to find the closest, according to some metric, possible successor state given the current state of the system. This however amounts to replace the dynamical ocean model by a purely data-driven model. A closely related approach was proposed by [30] based on the so-called Takens approach, replacing the dynamical model with a delay coordinate embedding model. Another technique, known as the Dynamic Ensemble Update (DEU), was introduced in the context of an EnKF [53], but uses a particular dictionary of sparse realizations to sparsify the filter estimate. Other approaches also resorted to some dictionaries to account for some missing physics [3], or to simplify the complexity of the dynamical model [10].

The approach we propose here is somehow different; it uses the full dynamical model for forecasting and the dictionary to provide a possible set of realizations (ensemble members) that represents the current uncertainties based on the forecast. We present and discuss different metrics to select the new ensemble members from the dictionary, and test their relevance with a realistic ensemble data assimilation exercise using a high resolution MITgcm model in the Red Sea. The paper is organized as follows. Section 2 recalls the EnKF and EnOI algorithms. Section 3, presents the adaptive EnOI algorithm and discusses approaches to select the ensemble members from the available dictionary. Section 4 presents the results of the implementation of the adaptive EnOI algorithm with the Lorenz 63 and 96 models. Section 5 describes the general circulation ocean model and its configuration, as well as the assimilated observations. It also outlines the design of the conducted

assimilation experiments, and discusses the filters performances and results. Finally, Section 6 concludes the work with a summary of the main findings and a discussion on the future directions.

2. Ensemble Data Assimilation

The data assimilation problem with the Ensemble Kalman filter is described following the state-space model formulation

$$\mathbf{x}_{t+1} = \mathbf{M}_t(\mathbf{x}_t) + \boldsymbol{\eta}_t, \quad (1)$$

$$\mathbf{y}_t = \mathbf{h}_t(\mathbf{x}_t) + \boldsymbol{\epsilon}_t, \quad (2)$$

where \mathbf{M}_t denotes the model representing the ocean dynamics, \mathbf{x}_t is the state vector at time t , and $\boldsymbol{\eta}_t$ is the model error. \mathbf{y}_t is the observation vector, which is related to the state via the measurement operator \mathbf{h}_t and $\boldsymbol{\epsilon}_t$ represents the observational error. Both $\boldsymbol{\eta}_t$ and $\boldsymbol{\epsilon}_t$ are assumed independent and normally distributed of mean zero and covariance matrices \mathbf{Q}_t and \mathbf{R}_t , respectively.

As a variant of the well-known Kalman filter (KF) [29], the EnKF represents the statistics (first two moments) of the system state by a collection of random realizations, or ensemble members [13, 21]. The estimate at any given time is then given by the sample mean and the error covariance is approximated by the sample covariance of the ensemble [13]. Here we adopt a deterministic formulation of the EnKF [21, 24]. Given a forecast ensemble of N members at time step t forming the matrix $\mathbf{X}_t^f = [\mathbf{x}_{1,t}^f, \dots, \mathbf{x}_{N,t}^f]$, with $\mathbf{x}_{i,t}^f$ denoting the i -th ensemble member at time t . The forecast ensemble anomaly is

$$\mathbf{X}_t^{f'} = \mathbf{X}_t^f - \frac{1}{N} \left(\sum_{i=1}^N \mathbf{x}_{i,t}^f \right) \mathbf{e}_{1 \times N}, \quad (3)$$

with $\mathbf{e}_{1 \times N}$ denoting the matrix with ones as elements and size $1 \times N$. At the analysis step, once an observation \mathbf{y}_t^o becomes available, the forecast state \mathbf{x}_t^f , which is the mean of \mathbf{X}_t^f , is updated using the standard Kalman filter correction step to obtain the analysis state

$$\mathbf{x}_t^a = \mathbf{x}_t^f + \mathbf{K}_t \left(\mathbf{y}_t^o - \mathbf{h}_t(\mathbf{x}_t^f) \right), \quad (4)$$

where $\mathbf{K}_t = \mathbf{P}_t^f \mathbf{H}_t^T \left(\mathbf{H}_t \mathbf{P}_t^f \mathbf{H}_t^T + \mathbf{R}_t \right)^{-1}$ is the Kalman gain. The forecast error covariance \mathbf{P}_t^f is estimated as $\frac{1}{N-1} \mathbf{X}_t^{f'} \mathbf{X}_t^{f'T}$, and the associated analysis

error covariance as

$$\mathbf{P}_t^a = [(\mathbf{P}_t^f)^{-1} + \mathbf{H}_t^T \mathbf{R}_t^{-1} \mathbf{H}_t]^ {-1}.$$

In the EAKF formulation, a matrix \mathbf{A}_t is introduced such that $\mathbf{P}_t^a = \mathbf{A}_t \mathbf{P}_t^f \mathbf{A}_t^T$. Based on a judicious choice of \mathbf{A}_t , an analysis ensemble is then resampled as $\mathbf{X}_t^a = \mathbf{A}_t \left(\mathbf{X}_t^f - \frac{1}{N} \left(\sum_{i=1}^N \mathbf{x}_{i,t}^f \right) \mathbf{e}_{1 \times N} \right) + \frac{1}{N} \left(\sum_{i=1}^N \mathbf{x}_{i,t}^a \right) \mathbf{e}_{1 \times N}$ in such a way to match the analysis \mathbf{x}_t^a and covariance \mathbf{P}_t^a before it is integrated by the model (1) to compute the next forecast [21], allowing for a dynamical update of the estimation error. A new assimilation cycle starts once the new observation becomes available.

EnOI is the optimal Interpolation (OI) variant of the EnKF ([14]) in which a pre-selected ensemble remains static during all assimilation cycles, with no feedback from the assimilation system to modify the forecast (background) covariance ([52]). In the EnOI forecast step, only the analysis state is integrated by the dynamical model for forecasting, before it gets updated again with the incoming observation based on the pre-selected ensemble ([14]). EnOI therefore leads to a drastic computational cost reduction (by almost a factor N) compared to an EnKF, and no resampling step is needed after the analysis. It further does not suffer from the typical ensemble collapse of the EnKFs, which often requires artificial inflation of its ensemble [1]. This makes EnOI a computationally very efficient approach for ensemble assimilation and was shown to be particularly robust in numerous ocean applications ([6, 43, 60, 63]).

3. Adaptive EnOI

The use of representative background covariances is critical for the performance of any data assimilation scheme, as these should describe the spatial and multivariate structure of the subspace in which the update with the observation is performed ([19, 34]). In particular, the behavior of ensemble assimilation methods largely depends on the representativeness of their (forecast) ensembles, based on which the background covariance is estimated ([56]). The ensemble should (i) describe the directions of estimation errors growth, and therefore be time-variant to follow their dynamical evolution ([21, 23]), and (ii) be large enough to infer reliable statistics between the observations and the forecast state and to provide enough rank (degrees of freedom) to fit the data ([19, 22]). In realistic large scale applications with general circulation ocean models, however, EnKFs can be only implemented with relatively limited ensembles O(10 members) to maintain a

manageable computational load ([21]). This usually results in rank-deficient background covariances that require various auxiliary techniques, such as covariance localization ([28]) inflation ([19]), to infer reasonable forecast increments from the incoming observations. Localization restricts the action of the increments only to grid points close to the observation, which helps increasing the background covariance rank and filters spurious correlations [27]. Another typical concern with ensemble data assimilation systems is the loss of spread in the forecast ensemble, which is associated with the dissipative nature of ocean models ([25, 32]) and the often misrepresented sources of model errors ([26]). This is often mitigated through simple ensemble inflation and/or stochastic perturbations (of the parameters and inputs) techniques [14].

EnOI schemes efficiently resolve the issue of computation load, which enables the use of large ensembles without cost increase. A static ensemble may however not always be representative of the modeled dynamics, especially when dealing with rapidly varying dynamics and those that experience sudden regime changes ([23]). To deal with the pronounced seasonality of the South China Sea, [63] suggested pre-selecting seasonal representative ensembles and then use these in an EnOI according to the season during which the observations are assimilated.

We propose here to push further the idea of using a time-varying ensemble in EnOI, not only by utilizing static ensembles by selecting a new ensemble on a seasonal basis, but at every assimilation cycle to account for the mesoscale and eventually intra-seasonal variability. We propose here to select the new ensemble after every forecast step from an available historical set, or “dictionary”, of ocean states describing the variability of the studied basin. The selection of the ensemble will be based on the best available information at the time of the update step, which in the context of an EnOI is the forecast state. The proposed assimilation workflow is schematized in Figure 1 and we will refer to it as the Adaptive EnOI, or AEnOI. The selection of the ensemble members from the dictionary is the corner stone of the proposed approach and is discussed in the next section.

3.1. Ensembles selection

We present two different strategies to select the ensemble members from an available dictionary: (i) select the elements that are the “closest” to the forecast according to a certain distance, (ii) select the elements that describe at best the filtering error subspace [42], based on the so-called Orthogonal Matching Pursuit (OMP) algorithm. After selecting the new members, the mean of the ensemble is replaced by the forecast state in both approaches,

so that only the ensemble anomaly is used in the EnOI algorithm. The incoming observations are not used in the selection, so that the data are not involved in the choice of the prior.

3.1.1. *Distance-based similarity selection*

We look for the dictionary elements that bear spatial similarities, or are the closest in some sense, to the forecast state according to a distance measure. The idea is that if an ocean state displays similar spatial features as the forecast state, it is also expected to carry information about the uncertainties around the forecast. One straightforward way to evaluate the distance between the forecast and the dictionary elements is to use the L2-norm, or L1-norm as illustrated in Figure 2. In our experiments, the assimilation results were quite close whether using L1-norm or L2-norm, and thus we only report here the results of the latter in the numerical experiments presented in Section 5.

Quantifying the similarity between two fields according to some norm may under-represent some localized ocean features in the overall basin-distance. We have also tried to involve correlations in our elements selection, but the strong environmental gradient in the Red Sea ([70]) dominated the correlations and made it difficult to distinguish the dictionary elements in this basin at the mesoscales.

3.1.2. *Error-subspace selection*

The basic idea is to identify a subset of the dictionary elements that represents at best the forecast error subspace in which the Kalman filter update is applied [23, 32, 41]. Here we propose to use a Matching Pursuit (MP) method, an interactive greedy algorithm that finds the *best matching* projections of a high-dimensional signal onto the span of a (complete) dictionary [37]. By selecting the elements that are most correlated with the current residuals (see Figure 3 for schematic illustration and algorithm’s description), MP attempts to approximately represent a signal using a sparse linear combination of the dictionary elements, called atoms, while minimizing the signal representational error in the dictionary. This is different than selecting the elements that are most correlated with the forecast state, and should lead to an ensemble with more spread describing the forecast state variability, assumingly representative of the filter error-subspace. In the Orthogonal Matching Pursuit (OMP) algorithm, the residual is always orthogonal to the span of the dictionary elements already selected. This can conceptually be implemented using a Gram-Schmidt scheme and results in convergence for a n -dimensional vector after at most p -iterations ($p \leq n$,

being the sparsity level) [61]. Enforcing orthogonal elements helps to avoid selecting redundant elements and provides more ensemble spread [61].

3.2. Implementation of the ensemble selection strategies

All the selection methods share the same workflow and the difference appears only at the selection stage (3.1. of Algorithm 1). The generic form of the ensemble selection is detailed in Algorithm 2, while Table 1 outlines two implementations of Algorithm 2, one with the L2 selection method and the other one with the OMP.

4. Preliminary experimentation with Lorenz models

The adaptive EnOI schemes are first tested and compared with the standard EnOI and EnKF algorithms using Lorenz-63 [35] (hereafter L-63), and Lorenz-96 [36] (hereafter L-96) models, two popular prototypes for assessing new assimilation schemes. For each model, all EnOI-based schemes use the same dictionary, constructed from a collection of samples that came with the EnKF-Matlab software [51], and was generated by a long model run. We conduct twin-experiments where the trajectory of a reference run is taken as the “true” trajectory from which synthetic observations are generated by adding zero-mean Gaussian white noise with variance σ_y^2 . The filters’ performances are evaluated using the root-mean-square error (RMSE) between the reference states and the filters’ estimates averaged over all variables and over the whole assimilation period. We implement all filters using the covariance inflation [2]. We further apply a local analysis [27] in the L-96 experiments.

4.1. Numerical experiments with L-63

The L-63 model is a nonlinear dissipative dynamical system that mimics an atmospheric chaotic behavior [45]. It is governed by the following systems of differential equations

$$\begin{cases} \frac{dx}{dt} = \sigma(y - x), \\ \frac{dy}{dt} = (\rho - z)x - y, \\ \frac{dz}{dt} = xy - \beta z, \end{cases}$$

where $\sigma = 10$, $\rho = 28$, and $\beta = 8/3$. The state variables x , y and z measure, respectively, the intensity of convective motion, the temperature difference between the ascending and descending currents, and the distortion of vertical temperature profile from linearity. The model is integrated using a fourth-order Runge-Kutta integration scheme, with a time step of 0.01 time units. After a spin-up period of roughly 20 days to remove any detrimental impact, the simulations are run for a period of five years in model time (i.e., 36500 model steps). We consider the case where all three variables are observed with $\sigma_y^2 = 2$. All schemes are tested with different values of inflation, ranging between 1 and 1.3, and compared, based on their minimum RMSEs. Our numerical experiments suggested that EnOI-based schemes do not require inflation, which was therefore not applied.

Figure 4 plots time series of the analysis RMSEs (upper subplot) and the forecast ensemble standard deviations (bottom subplot), as resulting from all schemes between assimilation steps 3000 and 3500, for illustration. Data are assimilated every 4 model steps and the ensemble size is set to 100. The RMSEs time series suggest a clear outperformance of the EnKF, followed by the AEnOI-L2 and the AEnOI-OMP. The reported RMSE values, averaged over the whole simulation period (0.172 for EnKF, 1.032 for AEnOI-L2, 1.119 for AEnOI-OMP, and 1.205 for EnOI), further confirm this and support our expectations about the capabilities of the AEnOI schemes in improving the EnOI performances, although they all fall behind EnKF. This is expected as the EnKF evolves the underlying state distribution by updating the ensemble members with the model dynamics. This however might be computationally demanding, since large ensembles are usually needed to properly describe the state statistics. For example, a one-year filtering run with 200 members, assimilating every 50 model steps (2.5 days), completed in 6.5538 s with the EnKF, 0.5103 s with the AEnOI-L2 and 2.9924 s with the AEnOI-OMP. Regarding the ensemble spread, the results suggest that EnKF, and to a lesser extent AEnOI-L2, exhibit the smallest spreads. EnOI and AEnOI-OMP however have larger spreads, but with different patterns. Indeed, EnOI has of course a constant spread, whereas AEnOI-OMP suggests a strongly variable spread over time. We further study the schemes' sensitivities to different ensemble sizes (N_e) and frequencies of assimilation in Figures 5 and 6, respectively. Overall, increasing the ensemble size reduces the RMSE values for all schemes, except for AEnOI-L2, which yields the most accurate results with small ensembles (10, 20 and 40), while for these ensemble sizes, AEnOI-OMP does not improve EnOI performances. As the ensemble size increases, the benefit from AEnOI-OMP becomes clearer and its performances approach those of AEnOI-L2 while both schemes remain

better than EnOI. Similarly, assimilating the data more frequently further improves the results although the EnOI-based schemes seem less sensitive to the assimilation period than the EnKF.

4.2. Numerical experiments with L-96

The L-96 model simulates the time evolution of an atmospheric quantity based on a set of differential equations:

$$\frac{dx_k}{dt} = (x_{k+1} - x_{k-2})x_{k-1} - x_k + F, \quad k = 1, \dots, K. \quad (5)$$

where x_k denotes the k^{th} element of the state \mathbf{x} . The nonlinear (quadratic) terms represent advection and the linear term simulates dissipation. In its most common form, the system dimension is $K = 40$ and the forcing term F is set to 8, a value for which the model exhibits a chaotic behavior. Boundary conditions are periodic (i.e; $x_{-1} = x_{39}$, $x_0 = x_{40}$ and $x_{41} = x_1$). The model is integrated using a fourth-order Runge-Kutta integration scheme, with a time step of 0.05 time units. After a spin-up period of roughly 20 days to remove any detrimental impact, the simulations are run for a period of five years in model time (i.e., 7300 model steps). We consider the case where all variables are observed with $\sigma_y^2 = 1$. We test the schemes using different values of inflation ranging between 1 and 1.3. We apply the standard local analysis approach by restricting the update of each grid point to only observations falling within some influence radius ([50]). The localization support radii vary from 2 (strong localization) to 40 (weak localization) grid points. The schemes are then compared based on their minimum RMSEs over all possible combinations of inflation and localization values. Figure 7 gives an idea about the schemes' needs for inflation and localization by plotting the RMSEs as a function of the localization radius and inflation factor.

Based on an extensive set of numerical experiments, the EnKF was the most sensitive to the choice of inflation and localization. Regarding EnOI-based schemes, the results suggest that they have enough spread and therefore do not require inflation. AEnOI-OMP was the least dependent on localization among the EnOI schemes. A sensitivity further analysis suggests that increasing the ensemble size and the assimilation frequency enhances the schemes behaviors (Figures 8,9), with less sensitivity from the EnOI-based schemes. The results further suggest that, in all tested scenarios, EnKF provides the most accurate estimates. AEnOI-OMP yields similar results as EnOI whereas AEnOI-L2 clearly improves its performances. One may also notice that, when data are assimilated less frequently using a small

ensemble, the performance of EnKF degrades and approaches those of the EnOI-based schemes. Therefore, in these challenging cases, the benefit from using AEnOI-L2 as a computationally less demanding alternative to EnKF, becomes more pronounced.

5. Experimentation with an ocean general circulation model in the Red Sea

5.1. The ocean model

We employ the Massachusetts Institute of Technology general circulation model (MITgcm), which solves the Navier-Stokes equations under the implicit free surface and Boussinesq approximations ([38]). The model is configured for the domain 30°E - 50°E and 10°N - 30°N covering the whole Red Sea, including the Gulf of Suez, the Gulf of Aqaba, and part of the Gulf of Aden where an open boundary connects it to the Arabian Sea. The model is configured with a Cartesian grid at an eddy-resolving grid resolution of $0.04^{\circ} \times 0.04^{\circ}$. In the vertical direction, the configuration has 50 layers, with 4 m spacing in the surface and 300 m near the bottom. The bathymetry of the basin is derived from the General Bathymetric Chart of the Ocean (GEBCO, available at http://www.gebco.net/data_and_products/gridded_bathymetry_data). The model uses a direct space time 3^{rd} order scheme for tracer advection, harmonic viscosity with coefficients of $30 \text{ m}^2/\text{s}$ in the horizontal and $7 \times 10^{-4} \text{ m}^2/\text{s}$ in the vertical, implicit horizontal diffusion for both temperature and salinity, and the K-Profile Parameterization (KPP) scheme ([31]) for vertical mixing with a vertical diffusion coefficient of $10^{-5} \text{ m}^2/\text{s}$ for both temperature and salinity. The open boundary conditions for temperature, salinity, and horizontal velocity are prescribed daily from the Global Ocean Reanalysis and Simulation data (GLORYS; [47]) available on a $1/12^{\circ}$ horizontal grid. A sponge layer of 5 grid boxes with a relaxation period of 1-day is implemented for smooth incorporation of open ocean conditions through the eastern boundary. The normal velocities at the boundary are adjusted to match the volume flux of GLORYS, which is estimated from GLORYS sea surface height (SSH) variations inside the model domain. The resulting inflow at the eastern boundary ensures consistency between the model and GLORYS basin-scale SSH. The model was spun-up for 31 years starting from 1979 to 2010 using the European Center for Medium Range Weather Forecast (ECMWF) reanalysis of atmospheric surface fluxes of radiation, momentum, freshwater sampled every 6-hour and available on a $75 \text{ km} \times 75 \text{ km}$ grid ([7]). The model has been extensively validated for the Red Sea by earlier studies

(e.g. [16, 60, 65, 66, 67]). For comparison with the assimilation runs (as further discussed in the next section), the same model configuration was integrated forward for the year 2011 using 6-hourly ECMWF atmospheric fields available at $50 \text{ km} \times 50 \text{ km}$ resolution. We refer to this model free-run experiment without assimilation as *Fexp*.

5.2. Experimental setup

Available observations are assimilated using the Ensemble Adjustment Kalman Filter (EAKF) available in the DART-MITgcm (Data Assimilation Research Testbed) package ([20, 24]) implemented for the Red Sea by [60]. All the experiments, in the present study, assimilate the data every 3 days, using a $\sim 300 \text{ km}$ horizontal localization radius and a multiplicative inflation factor of 1.1, as suggested by [60]. We assimilated observations of SST data generated from a level-4 daily $0.25^\circ \times 0.25^\circ$ resolution product of [49] (which was prepared by blending SST measurements from in situ and advanced very high resolution radiometer infrared satellites), and along-track satellite level-3 merged altimeter filtered sea level anomalies (SLA; corrected for dynamic atmospheric, ocean tide, and long wavelength errors) from Copernicus Marine Environment Monitoring Service (CMEMS; [39]). To compute the innovations between the SLA observations and the model SSH during assimilation, we add the model mean SSH to SLA observations prior to assimilation. Observations errors are assumed uncorrelated, and are prescribed with error variance of $(0.04 \text{ m})^2$ for SLA, and vary between $(0.1 \text{ }^\circ\text{C})^2$ and $(0.6 \text{ }^\circ\text{C})^2$ for SST in accordance with the interpolation errors specified in the level-4 gridded SST product of [49]. Four different assimilation experiments were conducted under the same conditions: EAKF with 50 members, and EnOI, AEnOI-L2, and AEnOI-OMP with 300 ensemble members. They differ only in terms of the underlying method to sample/select the ensemble from a long dictionary of MITgcm outputs simulated during the period 2002-2016. EAKF dynamically evolves the ocean ensemble. Its initial ensemble is generated by first selecting *Fexp* fields corresponding to ± 15 days from January 1st and then by adjusting the ensemble mean to the same initial state as *Fexp*. The EnOI uses a static ensemble of 300 members across all assimilation cycles (60 cycles in total) by selecting ocean states of 2002-2016 model hindcasts. AEnOI-L2 and AEnOI-OMP, dynamically select 300 members, based on the SST distance between the current ocean state and the dictionary elements. This choice is motivated by two factors: the SST exhibits a seasonal signal, and the ensemble members selection would have been computationally demanding (especially for OMP) if based on the full ocean state vector (10^7) and a dictionary with large number of

elements. All the EnOI assimilation experiments are conducted over a 6-month period in 2011, starting from January 1st, 2011 using the same initial condition as *Fexp*.

Unless stated, we analyze daily averaged forecasts as they result from the different assimilation experiments. Bias, correlations and root-mean-square-errors (RMSE) of the assimilated solution (both analysis snapshots and daily averaged forecasts) for SST and SSH are computed with respect to the merged satellite level-3 observations of the Group for High Resolution Sea Surface Temperature (GHRSSST; [12]) and merged along track level-3 altimeter observations of SSH from CMEMS ([39]), respectively. In order to demonstrate the relative performance of the assimilation system with respect to interpolated products, we employed level-4 SST and SSH products. The interpolated SST product is a high-resolution daily averaged level-4 SST product from OSTIA (Operational Sea Surface Temperature and Sea Ice Analysis) [8, 57], generated on a 0.054° (~ 6 km) grid by combining SST data from various satellites and in situ observations using an Optimal Interpolation (OI) system. The interpolated SSH product is the multi-mission altimeter merged satellite level-4 gridded Absolute Dynamic Topography (ADT) provided by CMEMS (here after CMEMS-L4; [39]), which is also available daily at a resolution of $0.25^\circ \times 0.25^\circ$. The maximum reported ADT mapping error (provided along with the CMEMS-L4 ADT product) during the analysis period 1st January-30th June, 2011 is estimated between 1.8 cm - 4 cm in the southern Red Sea and reaches up to 7 cm in the northern Red Sea. In order to use it in the present study, we adjust the CMEMS-L4 ADT by replacing its 15-year average by the model equivalent sea surface height climatology, similarly to the treatment of the level-3 SLA observations for assimilation.

5.3. Experimental results

Figure 10 displays spatial maps of SST forecast statistics (computed over the study period, i.e. 1st January to 30th June, 2011), compared to satellite level-3 SST observations, from different assimilation schemes, the free model experiment and the interpolated data product. The results indicate that the standard deviation (STD) is large over the Gulf of Aden (reaching 2 °C) and small over the central parts of the Red Sea (below 1.2 °C). Outputs from the free model run captures this contrasting feature. However, it underestimates the SST STD over the whole domain, particularly in the northern and central parts of the Red Sea. Those underestimations of SST STD, as well as SST biases, are improved by assimilation. The improvements are

more pronounced in the adaptive and ensemble optimal interpolation experiments relative to the EAKF. However, the EnOI and AEnOI-OMP suggest increased SST biases in the Gulf of Aden. Root-mean-square errors (RMSEs) and correlations also deteriorated, particularly in EnOI, with RMSEs increasing from 0.5 °C to 1 °C and correlations dropping from 0.95 to 0.8. Assimilation with the AEnOI-L2 strategy, on the other hand, yields SST improvements, with biases and RMSEs mostly within 0.5 °C and correlations above 0.8, all over the model domain, including Gulf of Aden and the Red Sea. AEnOI-L2 results are even better than the interpolated SST product, particularly in the northern and central Red Sea.

We further analyzed the time evolution of RMSEs for the daily averaged SST forecasts (Figure 11a) and for 3-day spaced SST analysis snapshots (Figure 11b) corresponding to the studied domain. As shown in Figure 11a, RMSEs of SST forecasts from all the model experiments and interpolation products does exhibit time dependence, with SST RMSEs dipping in February and peaking during June, except for EnOI and AEnOI-OMP which showed an additional peak (reaching 2 °C and 1.6 °C in EnOI and AEnOI-OMP, respectively) during the month of March. Interestingly, SST RMSEs resulting from EnOI and AEnOI-OMP are larger than those of *Fexp* until the last week of April. SST RMSEs are almost always less than those of *Fexp* when assimilating observations with EAKF, but they are further improved, even over the interpolated product, with the AEnOI-L2 strategy. Assimilation fits the observations better in AEnOI-L2 than in EAKF (Figure 11b), which seem to be due to improved SST spread (as discussed in the subsequent paragraphs using Figure 13), explaining the better SST forecasts in AEnOI-L2. The SST analyses of all the three ensemble optimal interpolation experiments are indeed almost identical (Figure 11b). The failure to yield uniformly low SST forecast RMSEs and the occasional SST degradations in EnOI and AEnOI-OMP compared to the consistent improvements witnessed in AEnOI-L2 and EAKF may be attributed to the repercussion from comparatively larger dynamical imbalances in EnOI and AEnOI-OMP analyses (as discussed in the subsequent paragraphs) [33, 55].

Figures 11c and 11d, respectively, display the time evolution of SSH RMSEs for daily averaged forecasts and 3-day spaced analysis snapshots from different experiments and interpolated product. SSH RMSEs of *Fexp* exhibit noticeable fluctuations with largest values (reaching 14 cm) during January and smallest values (~ 5 cm) during the end of May. Unlike the free model, SSH RMSEs in the interpolated product are stable with values around 5 cm. Assimilating observations with EAKF, EnOI or AEnOI-OMP also yields SSH RMSEs close to 5 cm, but they exhibit fluctuations in SSH

RMSEs although not as large as *Fexp*. The fluctuations are reduced in AEnOI-L2, and the SSH RMSEs are generally lower than those of the interpolated product. In order to spatially investigate the assimilation results for SSH we analyzed the region wise statistics. Since the altimeter coverage is too sparse over the model domain to yield spatial maps of statistics, we tabulated (Table 2) statistics for four different regions: Gulf of Aden (GoA; 30°E-50°E and 10°N-14°N), Southern Red Sea (SRS; 30°E-50°E and 14°N-19°N), Central Red Sea (CRS; 30°E-50°E and 19°N-23°N) and Northern Red Sea (NRS; 30°E-50°E and 23°N-28°N). *Fexp* underestimates the STD (up to 3 cm) and the mean (up to 8 cm) of SSH. The underestimations of the mean are largest in the NRS (by 160%) and the largest underestimations of the STD are in the SRS (by 27%). SSH RMSEs (9-11 cm) and correlations (0.4-0.86) are also poor in *Fexp*, particularly in the GoA and the NRS. The interpolated SSH product also underestimates the mean, but provides robust estimates of the STD, with low RMSEs (5-6cm) and high correlations (0.94-0.98) throughout the domain. Assimilation improves the SSH mean and STDs considerably throughout the domain, even better than (or on par with) the interpolated data product. SSH RMSEs (5-7 cm) and correlations (0.54-0.92) are also improved compared to *Fexp*, and still less than the interpolated product. Interestingly, AEnOI-OMP (EAKF) improvements are less pronounced than those resulting from the standard EnOI, which is probably related to the SSH spread of the background ensemble, as further discussed in the subsequent paragraphs. The differences between EnOI and AEnOI-L2 are not so large except for the GoA, in which AEnOI-L2 yields better results than the rest of the assimilation schemes.

We also examine the estimated ocean state in the subsurface to assess the impact of the assimilation strategies in these sparsely observed layers. The ocean state in the subsurface layers is noisy in EnOI (Figures 12e and 12f) compared to *Fexp* (Figures 12a and 12b) and to EAKF (Figures 12c and 12d), consistent with the results of [60], in which the noise in the subsurface was attributed to pronounced dynamical imbalances in the analysis. While the ocean state becomes noisier in AEnOI-OMP (Figures 12i and 12j), AEnOI-L2 (Figures 12g and 12h) reduces this noise and yields more organized subsurface structures. For instance, EnOI simulates abrupt jumps in the 22 °C isotherm in the months of March, April, and also in May (Figure 12k) at (38°E, 22°N), and these are more frequent and larger in AEnOI-OMP. Such abrupt jumps do not appear in the results of AEnOI-L2, indicating a more stable solution. Dynamical imbalances (noise) may result from inappropriate analysis update due to spurious spread, and correlations in the background ensemble. These aspects are further discussed

in the next paragraphs.

Figure 13 plots the spatial distribution of the ensemble spread on 1-May-2011 from the different filtering schemes. The ensemble spreads of SSH, SST and subsurface temperature are considerably larger in all the ensemble optimal interpolation assimilation experiments compared to those of EAKF. This is because the spread introduced in the ensemble of initial conditions in EAKF fades out after few analysis cycles, and because the ensemble optimal interpolation strategies do not lose spread as they select members from model hindcasts after each analysis cycle. The spreads of SSH, SST and upper layer temperatures resulting from EnOI, AEnOI-L2 and AEnOI-OMP are significantly different. AEnOI-L2 selects the ensemble members from a broader range of months based on their closeness to the forecast SST, which seem to result here in a small ensemble spread (Figure 13g and 13k). AEnOI-OMP selects the ensemble members based on the correlations of the dictionary elements with residuals of the forecast state in the ensemble subspace (weaker the correlation better are the chances for selection). As a result, the selected members are not necessarily correlated/close to the forecasted SST, and may thus exhibit larger ensemble spread (Figures 13d, 13h and 13l). Large ensemble spreads may cause a data overfit, and amplify the noise in the filter updates, particularly in the data sparse regions ([44, 54]). This may explain the more (less) abrupt jumps in the 22 °C isotherm in AEnOI-OMP (AEnOI-L2) compared to EnOI.

One of the key assumptions of an EnKF framework lies on Gaussian forecast errors, based on which the members are updated with the observations using the Kalman linear analysis step ([25]). In the ensemble optimal interpolation schemes, the forecast error is estimated based on the anomalies of the selected ensemble. We assess the relevance of the Gaussian assumption in our setup by analyzing the histogram of SST ensembles at three locations in the northern, central and southern Red Sea on 1-May-2011, as shown in Figure 14. At all these locations, the prior distributions in EnOI and AEnOI-L2 are clearly more Gaussian than that of the AEnOI-OMP. The OMP scheme provides a more scattered ensemble that is far from a Gaussian distribution, and this may limit the relevance of the Kalman-based update step.

We also analyzed the SST correlation range at three different locations in the northern, central and southern Red Sea for the different EnOI schemes (Figure 15). At all locations, the SST correlation range for AEnOI-L2 is narrower and less noisy than those of EnOI, and AEnOI-OMP, suggesting less spurious long-range correlations. This means that AEnOI-L2 could be configured with a larger localization radius, which may subsequently result

in more dynamically consistent ocean state estimates ([15]). Given that AEnOI-L2 only forecasts the analysis state, this would enable using larger ensembles to rely even less on localization ([59]), without significantly increasing the computational cost. In our specific system, one MITgcm model run requires 4.8 core hours for a 3-day simulation and a DART-filter update requires 111 core hours. Therefore, one EAKF assimilation step with 300 (50) members consumes 1551 (351) core hours. The adaptive schemes involve a single model run for forecasting, and the selection step of its 300-member ensemble requires 21.37 and 20.77 core hours for the AEnOI-L2 and AEnOI-OMP, respectively, followed by a filter update. This amounts to an approximate computational cost of 137 core hours for each of the adaptive schemes and translates to more than a factor 10 (2) cost saving compared to the EAKF.

6. Conclusions

The Red Sea is characterized by a marked seasonal variability and strong mesoscales activity. In order to account for these variations at different time scales with reasonable computational burden, we proposed new cost-effective adaptive Ensemble Optimal Interpolation (AEnOI) schemes for assimilating multivariate data sets of the Red Sea based on the Data Assimilation Research Testbed (DART) and the MIT general circulation model (MITgcm).

The AEnOI schemes select the ensemble members from a complete dictionary describing the underlying system variability. The members selection is based on their similarity to, according to a certain criteria, or to their representativeness of the current forecast state, which represents the best available information at the time of the incoming observations. Two approaches for selecting the ensemble members were proposed: the first is based on the L2-distance between the forecast and the dictionary elements, and the second uses an Orthogonal Matching Pursuit (OMP) algorithm to identify the error-subspace of the forecast state. In term of computational efficiency, EnOI was of course an advantage since the selection process is applied offline and only once, before the start of the assimilation experiments. The AEnOI schemes enable however for adaptive selection of the ensemble members, which could account for instance for inter-seasonal and mesoscale variability.

The AEnOI schemes were first implemented and validated with the Lorenz-63 (L-63) and the Lorenz-96 (L-96) models, compared against the Ensemble Kalman filter (EnKF) and the standard EnOI. While the EnKF

yields the best results, eventually at the expense of applying auxiliary techniques such as inflation and localization, and higher computational cost, AEnOIs generally yield more accurate estimates than the standard EnOI, in terms of RMSE. They are further, particularly AEnOI-L2, computationally very efficient and may provide an alternative to the EnKF in the challenging scenario of small ensembles.

Within the DART-MITgcm Red Sea assimilation system, the AEnOI schemes operate on a dictionary of ocean realizations describing the multiscale temporal and spatial variability of the basin. Different aspects of the assimilation system have been assessed; including SST and SSH biases, standard deviations, correlations, and root-mean-square errors. AEnOI-L2 yields substantial improvements in certain regions of the Red Sea, whereas the AEnOI-OMP and the EnOI lead, in general, to more or less comparable assimilation results in our particular domain.

The AEnOI schemes, AEnOI-L2 more precisely, provided competitive performances to the computationally much demanding ensemble (adjustment) Kalman filter, especially in situations when the model forward integration is computationally demanding. We will work in the future on developing adaptive Hybrid schemes in which a new ensemble member will be selected from a dictionary, eventually regionally, based on the statistics of an (small) evolving ensemble. The resulting ensemble will combine the spread benefit of the EnOI scheme and will constrain it by that of the evolving ensemble that accounts for the error-of-the-day. We are also planning to implement these schemes within a stochastic EnKF framework based on the scheme proposed by [24].

7. Acknowledgment

This work was funded by the Office of Sponsored Research (OSR) at King Abdullah University of Science and Technology (KAUST) under the Virtual Red Sea Initiative (Grant #REP/1/3268-01-01) and the KAUST Center for Marine Environmental Observations (SAKMEO). The research made use of the KAUST supercomputing facilities.

- [1] Anderson, J. L., 2001. An ensemble adjustment kalman filter for data assimilation. *Monthly Weather Review* 129 (12), 2884–2903.
URL [https://doi.org/10.1175/1520-0493\(2001\)129<2884:AEAKFF>2.0.CO;2](https://doi.org/10.1175/1520-0493(2001)129<2884:AEAKFF>2.0.CO;2)
- [2] Anderson, J. L., Anderson, S. L., 1999. A monte carlo implementation of the nonlinear filtering problem to produce ensemble assimilations and forecasts. *Monthly Weather Review* 127 (12), 2741–2758.
URL [https://doi.org/10.1175/1520-0493\(1999\)127<2741:AMCIOT>2.0.CO;2](https://doi.org/10.1175/1520-0493(1999)127<2741:AMCIOT>2.0.CO;2)
- [3] Berry, T., Harlim, J., 2016. Forecasting turbulent modes with non-parametric diffusion models: Learning from noisy data. *Physica D: Nonlinear Phenomena* 320, 57 – 76.
URL <http://www.sciencedirect.com/science/article/pii/S0167278916000166>
- [4] Carvalho, S., Aylagas, E., Villalobos, R., Kattan, Y., Berumen, M., Pearman, J. K., 2019. Beyond the visual: using metabarcoding to characterize the hidden reef cryptobiome. *Proceedings of the Royal Society B: Biological Sciences* 286 (1896), 20182697.
URL <https://royalsocietypublishing.org/doi/abs/10.1098/rspb.2018.2697>
- [5] Clifford, M., Horton, C., Schmitz, J., Kantha, L. H., 1997. An oceanographic nowcast/forecast system for the red sea. *Journal of Geophysical Research: Oceans* 102 (C11), 25101–25122.
URL <http://dx.doi.org/10.1029/97JC01919>
- [6] Counillon, F., Bertino, L., 2009. Ensemble optimal interpolation: multivariate properties in the gulf of mexico. *Tellus A* 61 (2), 296–308.
URL <http://dx.doi.org/10.1111/j.1600-0870.2008.00383.x>
- [7] Dee, D. P., Uppala, S. M., Simmons, A. J., Berrisford, P., Poli, P., Kobayashi, S., Andrae, U., Balmaseda, M. A., Balsamo, G., Bauer, P., Bechtold, P., Beljaars, A. C. M., van de Berg, L., Bidlot, J., Bormann, N., Delsol, C., Dragani, R., Fuentes, M., Geer, A. J., Haimberger, L., Healy, S. B., Hersbach, H., Hólm, E. V., Isaksen, L., Kållberg, P., Köhler, M., Matricardi, M., McNally, A. P., Monge-Sanz, B. M., Morcrette, J.-J., Park, B.-K., Peubey, C., de Rosnay, P., Tavolato, C., Thépaut, J.-N., Vitart, F., 2011. The era-interim reanalysis: configuration and performance of the data assimilation system. *Quarterly*

Journal of the Royal Meteorological Society 137 (656), 553–597.
URL <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.828>

- [8] Donlon, C. J., Martin, M., Stark, J., Roberts-Jones, J., Fiedler, E., Wimmer, W., 2012. The operational sea surface temperature and sea ice analysis (ostia) system. Remote Sensing of Environment 116, 140 – 158, advanced Along Track Scanning Radiometer(AATSR) Special Issue.
URL <http://www.sciencedirect.com/science/article/pii/S0034425711002197>
- [9] Dreano, D., Raitos, D. E., Gittings, J., Krokos, G., Hoteit, I., Dec 2016. The gulf of aden intermediate water intrusion regulates the southern red sea summer phytoplankton blooms. PloS one 11 (12), e0168440–e0168440.
URL <https://www.ncbi.nlm.nih.gov/pubmed/28006006>
- [10] Dreano, D., Tsiaras, K., Triantafyllou, G., Hoteit, I., Jul 2017. Efficient ensemble forecasting of marine ecology with clustered 1d models and statistical lateral exchange: application to the red sea. Ocean Dynamics 67 (7), 935–947.
URL <https://doi.org/10.1007/s10236-017-1065-0>
- [11] Edwards, C. A., Moore, A. M., Hoteit, I., Cornuelle, B. D., 2015. Regional ocean data assimilation. Annual Review of Marine Science 7 (1), 21–42, pMID: 25103331.
URL <https://doi.org/10.1146/annurev-marine-010814-015821>
- [12] EUMETSAT/OSI-SAF, 2008. GHR SST Level 3P Global Subskin Sea Surface Temperature from the Advanced Very High Resolution Radiometer (AVHRR) on the MetOp-A satellite. <https://doi.org/10.5067/GHGMT-3PE01>, Ver. 1. PO.DAAC, CA, USA. Dataset accessed 2018.10.01.
- [13] Evensen, G., 1994. Sequential data assimilation with a nonlinear quasi-geostrophic model using monte carlo methods to forecast error statistics. Journal of Geophysical Research: Oceans 99 (C5), 10143–10162.
URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/94JC00572>

- [14] Evensen, G., Nov. 2003. The Ensemble Kalman Filter: Theoretical formulation and practical implementation. *Ocean Dynamics* 53 (4), 343–367.
URL <http://link.springer.com/10.1007/s10236-003-0036-9>
- [15] Flowerdew, J., 2015. Towards a theory of optimal localisation. *Tellus A: Dynamic Meteorology and Oceanography* 67 (1), 25257.
URL <https://doi.org/10.3402/tellusa.v67.25257>
- [16] Gittings, J. A., Raitsos, D. E., Krokos, G., Hoteit, I., 2018. Impacts of warming on phytoplankton abundance and phenology in a typical tropical marine ecosystem. *Scientific Reports* 8 (1), 2240.
URL <https://doi.org/10.1038/s41598-018-20560-5>
- [17] Guo, D., Akylas, T. R., Zhan, P., Kartadikaria, A., Hoteit, I., 2016. On the generation and evolution of internal solitary waves in the southern red sea. *Journal of Geophysical Research: Oceans* 121 (12), 8566–8584.
URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2016JC012221>
- [18] Guo, D., Kartadikaria, A., Zhan, P., Xie, J., Li, M., Hoteit, I., 2018. Baroclinic tides simulation in the red sea: Comparison to observations and basic characteristics. *Journal of Geophysical Research: Oceans* 123 (12), 9389–9404.
URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018JC013970>
- [19] Hamill, T. M., Snyder, C., 2000. A hybrid ensemble kalman filter3d variational analysis scheme. *Monthly Weather Review* 128 (8), 2905–2919.
URL [https://doi.org/10.1175/1520-0493\(2000\)128<2905:AHEKFV>2.0.CO;2](https://doi.org/10.1175/1520-0493(2000)128<2905:AHEKFV>2.0.CO;2)
- [20] Hoteit, I., Hoar, T., Gopalakrishnan, G., Collins, N., Anderson, J., Cornuelle, B., Köhl, A., Heimbach, P., 2013. A mitgcm/dart ensemble analysis and prediction system with application to the gulf of mexico. *Dynamics of Atmospheres and Oceans* 63, 1 – 23.
URL <http://www.sciencedirect.com/science/article/pii/S0377026513000249>
- [21] Hoteit, I., Luo, X., Bocquet, M., Köhl, A., Ait-El-Fquih, B., 2018. Data Assimilation in Oceanography: Current Status and New Directions.

GODAE OceanView, Ch. 17, pp. 465–512.
URL <https://doi.org/10.17125/gov2018.ch17>

- [22] Hoteit, I., Luo, X., Pham, D.-T., 2012. Particle kalman filtering: A nonlinear bayesian framework for ensemble kalman filters. *Monthly Weather Review* 140 (2), 528–542.
URL <https://doi.org/10.1175/2011MWR3640.1>
- [23] Hoteit, I., Pham, D.-T., Blum, J., 2002. A simplified reduced order kalman filtering and application to altimetric data assimilation in tropical pacific. *Journal of Marine Systems* 36 (1), 101 – 127.
URL <http://www.sciencedirect.com/science/article/pii/S092479630200129X>
- [24] Hoteit, I., Pham, D.-T., Gharamti, M. E., Luo, X., 2015. Mitigating observation perturbation sampling errors in the stochastic enkf. *Monthly Weather Review* 143 (7), 2918–2936.
URL <https://doi.org/10.1175/MWR-D-14-00088.1>
- [25] Hoteit, I., Pham, D.-T., Triantafyllou, G., Korres, G., 2008. A new approximate solution of the optimal nonlinear filter for data assimilation in meteorology and oceanography. *Monthly Weather Review* 136 (1), 317–334.
URL <https://doi.org/10.1175/2007MWR1927.1>
- [26] Hoteit, I., Triantafyllou, G., Korres, G., 04 2007. Using low-rank ensemble kalman filters for data assimilation with high dimensional imperfect models. *JNAIAM. Journal of Numerical Analysis, Industrial and Applied Mathematics* 2.
URL <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.418.6672>
- [27] Houtekamer, P. L., Mitchell, H. L., 1998. Data assimilation using an ensemble kalman filter technique. *Monthly Weather Review* 126 (3), 796–811.
URL [https://doi.org/10.1175/1520-0493\(1998\)126<0796:DAUAEK>2.0.CO;2](https://doi.org/10.1175/1520-0493(1998)126<0796:DAUAEK>2.0.CO;2)
- [28] Houtekamer, P. L., Mitchell, H. L., 2001. A sequential ensemble kalman filter for atmospheric data assimilation. *Monthly Weather Review* 129 (1), 123–137.
URL [https://doi.org/10.1175/1520-0493\(2001\)129<0123:ASEKFF>2.0.CO;2](https://doi.org/10.1175/1520-0493(2001)129<0123:ASEKFF>2.0.CO;2)

- [29] Kalman, R. E., 1960. A New Approach to Linear Filtering and Prediction Problems. *Journal of Basic Engineering* 82(1), 35–45.
URL <https://doi.org/10.1115/1.3662552>
- [30] Khaki, M., Hamilton, F., Forootan, E., Hoteit, I., Awange, J., Kuhn, M., 2018. Nonparametric data assimilation scheme for land hydrological applications. *Water Resources Research* 54 (7), 4946–4964.
URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018WR022854>
- [31] Large, W. G., McWilliams, J. C., Doney, S. C., 1994. Oceanic vertical mixing: A review and a model with a nonlocal boundary layer parameterization. *Reviews of Geophysics* 32 (4), 363–403.
URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/94RG01872>
- [32] Lermusiaux, P. F. J., Robinson, A. R., 1999. Data assimilation via error subspace statistical estimation. part i: Theory and schemes. *Monthly Weather Review* 127 (7), 1385–1407.
URL [https://doi.org/10.1175/1520-0493\(1999\)127<1385:DAVESS>2.0.CO;2](https://doi.org/10.1175/1520-0493(1999)127<1385:DAVESS>2.0.CO;2)
- [33] Li, Y., Toumi, R., 2017. A balanced kalman filter ocean data assimilation system with application to the south australian sea. *Ocean Modelling* 116, 159 – 172.
URL <http://www.sciencedirect.com/science/article/pii/S1463500317300963>
- [34] Lorenc, A. C., 2003. The potential of the ensemble kalman filter for nwp-a comparison with 4d-var. *Quarterly Journal of the Royal Meteorological Society* 129 (595), 3183–3203.
URL <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1256/qj.02.132>
- [35] Lorenz, E. N., 1963. Deterministic nonperiodic flow. *Journal of the Atmospheric Sciences* 20 (2), 130–141.
URL [https://doi.org/10.1175/1520-0469\(1963\)020<0130:DNF>2.0.CO;2](https://doi.org/10.1175/1520-0469(1963)020<0130:DNF>2.0.CO;2)
- [36] Lorenz, E. N., Emanuel, K. A., 1998. Optimal sites for supplementary weather observations: Simulation with a small model. *Journal of the Atmospheric Sciences* 55 (3), 399–414.

- URL [https://doi.org/10.1175/1520-0469\(1998\)055<0399:OSFSW0>2.0.CO;2](https://doi.org/10.1175/1520-0469(1998)055<0399:OSFSW0>2.0.CO;2)
- [37] Mallat, S. G., Zhang, Z., Dec 1993. Matching pursuits with time-frequency dictionaries. *IEEE Transactions on Signal Processing* 41 (12), 3397–3415.
URL <https://doi.org/10.1109/78.258082>
- [38] Marshall, J., Hill, C., Perelman, L., Adcroft, A., 1997. Hydrostatic, quasi-hydrostatic, and nonhydrostatic ocean modeling. *Journal of Geophysical Research: Oceans* 102 (C3), 5733–5752.
URL <https://doi.org/10.1029/96JC02776>
- [39] Mertz, F., Vinca, R., Caroline, M., Yannice, F., 2017. Product user manual, For sea level SLA products. Technical Report CMEMS-SL-PUM-008-032-051, issue 1.1.
URL <http://cmems-resources.cls.fr/documents/PUM/CMEMS-SL-PUM-008-032-051.pdf>
- [40] Nanninga, G. B., Saenz-Agudelo, P., Zhan, P., Hoteit, I., Berumen, M. L., Jun 2015. Not finding nemo: limited reef-scale retention in a coral reef fish. *Coral Reefs* 34 (2), 383–392.
URL <https://doi.org/10.1007/s00338-015-1266-2>
- [41] Nerger, L., Danilov, S., Kivman, G., Hiller, W., Schröter, J., 2004. Comparison of the ensemble kalman filter and the seik filter applied to a finite element model of the north atlantic. In: *EGU 1st General Assembly, Nice, France, April 25 - 30*.
URL <http://hdl.handle.net/10013/epic.20984>
- [42] Nerger, L., Schulte, S., Bunse-Gerstner, A., 2014. On the influence of model nonlinearity and localization on ensemble kalman smoothing. *Quarterly Journal of the Royal Meteorological Society* 140 (684), 2249–2259.
URL <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.2293>
- [43] Oke, P. R., Brassington, G. B., Griffin, D. A., Schiller, A., 2010. Ocean data assimilation: a case for ensemble optimal interpolation. *Australian Meteorological and Oceanographic Journal* 59, 67–76.
URL <https://doi.org/10.22499/2.5901.008>

- [44] Oke, P. R., Sakov, P., 2008. Representation error of oceanic observations for data assimilation. *Journal of Atmospheric and Oceanic Technology* 25 (6), 1004–1017.
URL <https://doi.org/10.1175/2007JTECH0558.1>
- [45] Palmer, T. N., 1993. Extended-range atmospheric prediction and the lorenz model. *Bulletin of the American Meteorological Society* 74 (1), 49–66.
URL [https://doi.org/10.1175/1520-0477\(1993\)074<0049:ERAPAT>2.0.CO;2](https://doi.org/10.1175/1520-0477(1993)074<0049:ERAPAT>2.0.CO;2)
- [46] Papadopoulos, V. P., Zhan, P., Sofianos, S. S., Raitzos, D. E., Qurban, M., Abualnaja, Y., Bower, A., Kontoyiannis, H., Pavlidou, A., Asharaf, T. T. M., Zarokanellos, N., Hoteit, I., 2015. Factors governing the deep ventilation of the red sea. *Journal of Geophysical Research: Oceans* 120 (11), 7493–7505.
URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2015JC010996>
- [47] Parent, L., Testut, C.-E., Brankart, J.-M., Verron, J., Brasseur, P., Gourdeau, L., 2003. Comparative assimilation of topeX/poseidon and ers altimeter data and of tao temperature data in the tropical pacific ocean during 19941998, and the mean sea-surface height issue. *Journal of Marine Systems* 40-41, 381 – 401, the Use of Data Assimilation in Coupled Hydrodynamic, Ecological and Bio-geo-chemical Models of the Ocean. Selected papers from the 33rd International Liege Colloquium on Ocean Dynamics, held in Liege, Belgium on May 7-11th, 2001.
URL <http://www.sciencedirect.com/science/article/pii/S0924796303000265>
- [48] Raitzos, D. E., Brewin, R. J. W., Zhan, P., Dreano, D., Pradhan, Y., Nanninga, G. B., Hoteit, I., 2017. Sensing coral reef connectivity pathways from space. *Scientific Reports* 7 (1), 9338.
URL <https://doi.org/10.1038/s41598-017-08729-w>
- [49] Reynolds, R. W., Smith, T. M., Liu, C., Chelton, D. B., Casey, K. S., Schlax, M. G., 2007. Daily high-resolution-blended analyses for sea surface temperature. *Journal of Climate* 20 (22), 5473–5496.
URL <https://doi.org/10.1175/2007JCLI1824.1>
- [50] Sakov, P., Bertino, L., Mar 2011. Relation between two common localisation methods for the enf. *Computational Geosciences* 15 (2), 225–

237.

URL <https://doi.org/10.1007/s10596-010-9202-6>

- [51] Sakov, P., Oliver, D. S., Bertino, L., 2012. An iterative enkf for strongly nonlinear systems. *Monthly Weather Review* 140 (6), 1988–2004.
URL <https://doi.org/10.1175/MWR-D-11-00176.1>
- [52] Sakov, P., Sandery, P. A., 2015. Comparison of enoi and enkf regional ocean reanalysis systems. *Ocean Modelling* 89, 45 – 60.
URL <http://www.sciencedirect.com/science/article/pii/S1463500315000219>
- [53] Sana, F., Katterbauer, K., Al-Naffouri, T., Hoteit, I., 2016. Orthogonal Matching Pursuit for Enhanced Recovery of Sparse Geological Structures with the Ensemble Kalman Filter. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 9 (4), 1710–1724.
URL <http://dx.doi.org/10.13140/RG.2.1.1803.8485>
- [54] Sanikommu, S., Banerjee, D. S., Baduru, B., Paul, B., Paul, A., Chakraborty, K., Hoteit, I., 2019. Impact of dynamical representational errors on an indian ocean ensemble data assimilation system. *Quarterly Journal of the Royal Meteorological Society* 145 (725), 3680–3691.
URL <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.3649>
- [55] Sivareddy, S., Paul, A., Sluka, T., Ravichandran, M., Kalnay, E., 2017. The pre-argo ocean reanalyses may be seriously affected by the spatial coverage of moored buoys. *Scientific Reports* 7 (1), 46685.
URL <https://doi.org/10.1038/srep46685>
- [56] Song, H., Hoteit, I., Cornuelle, B. D., Subramanian, A. C., 2010. An adaptive approach to mitigate background covariance limitations in the ensemble kalman filter. *Monthly Weather Review* 138 (7), 2825–2845.
URL <https://doi.org/10.1175/2010MWR2871.1>
- [57] Stark, J. D., Donlon, C. J., Martin, M. J., McCulloch, M. E., June 2007. Ostia : An operational, high resolution, real time, global sea surface temperature analysis system. In: *OCEANS 2007 - Europe*. pp. 1–4.
URL <https://doi.org/10.1109/OCEANSE.2007.4302251>
- [58] Tandeo, P., Ailliot, P., Ruiz, J., Hannart, A., Chapron, B., Cuzol, A., Monbet, V., Easton, R., Fablet, R., 2015. Combining Analog Method

and Ensemble Data Assimilation: Application to the Lorenz-63 Chaotic System. Springer International Publishing, Cham, pp. 3–12.
 URL http://dx.doi.org/10.1007/978-3-319-17220-0_1

- [59] Toye, H., Kortas, S., Zhan, P., Hoteit, I., 2018. A fault-tolerant hpc scheduler extension for large and operational ensemble data assimilation: Application to the red sea. *Journal of Computational Science* 27, 46 – 56.
 URL <http://www.sciencedirect.com/science/article/pii/S1877750317312905>
- [60] Toye, H., Zhan, P., Gopalakrishnan, G., Kartadikaria, A. R., Huang, H., Knio, O., Hoteit, I., Jul 2017. Ensemble data assimilation in the red sea: sensitivity to ensemble selection and atmospheric forcing. *Ocean Dynamics* 67 (7), 915–933.
 URL <https://doi.org/10.1007/s10236-017-1064-1>
- [61] Tropp, J. A., Gilbert, A. C., 2007. Signal recovery from random measurements via orthogonal matching pursuit. *Information Theory, IEEE Transactions on* 53 (12), 4655–4666.
 URL <https://doi.org/10.1109/TIT.2007.909108>
- [62] Vervatis, V., Testut, C., Mey, P. D., Ayoub, N., Chanut, J., Quattrocchi, G., 2016. Data assimilative twin-experiment in a high-resolution bay of biscay configuration: 4denoi based on stochastic modeling of the wind forcing. *Ocean Modelling* 100, 1 – 19.
 URL <http://www.sciencedirect.com/science/article/pii/S1463500316000044>
- [63] Xie, J., Zhu, J., 2010. Ensemble optimal interpolation schemes for assimilating Argo profiles into a hybrid coordinate ocean model. *Ocean Modelling* 33, 283–298.
 URL <https://doi.org/10.1016/j.ocemod.2010.03.002>
- [64] Yao, F., Hoteit, I., 2018. Rapid red sea deep water renewals caused by volcanic eruptions and the north atlantic oscillation. *Science Advances* 4 (6).
 URL <http://advances.sciencemag.org/content/4/6/eaar5637>
- [65] Yao, F., Hoteit, I., Pratt, L. J., Bower, A. S., Khl, A., Gopalakrishnan, G., Rivas, D., 2014. Seasonal overturning circulation in the red sea: 2. winter circulation. *Journal of Geophysical Research: Oceans* 119 (4),

2263–2289.

URL <http://dx.doi.org/10.1002/2013JC009331>

- [66] Yao, F., Hoteit, I., Pratt, L. J., Bower, A. S., Zhai, P., Khl, A., Gopalakrishnan, G., 2014. Seasonal overturning circulation in the red sea: 1. model validation and summer circulation. *Journal of Geophysical Research: Oceans* 119 (4), 2238–2262.
URL <http://dx.doi.org/10.1002/2013JC009004>
- [67] Zhan, P., Gopalakrishnan, G., Subramanian, A. C., Guo, D., Hoteit, I., 2018. Sensitivity studies of the red sea eddies using adjoint method. *Journal of Geophysical Research: Oceans* 123 (11), 8329–8345.
URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018JC014531>
- [68] Zhan, P., Krokos, G., Guo, D., Hoteit, I., 2019. Three-dimensional signature of the red sea eddies and eddy-induced transport. *Geophysical Research Letters* 0 (0).
URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018GL081387>
- [69] Zhan, P., Subramanian, A. C., Yao, F., Hoteit, I., 2014. Eddies in the red sea: A statistical and dynamical study. *Journal of Geophysical Research: Oceans* 119 (6), 3909–3925.
URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2013JC009563>
- [70] Zhan, P., Subramanian, A. C., Yao, F., Kartadikaria, A. R., Guo, D., Hoteit, I., 2016. The eddy kinetic energy budget in the red sea. *Journal of Geophysical Research: Oceans* 121 (7), 4732–4747.
URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2015JC011589>

Algorithm 1 Dictionary based schemes Data Assimilation Algorithm

0. Initialization: initial ensemble \mathbf{X}^f

1. Analysis step:

Input: \mathbf{X}^f

Output: \mathbf{x}^a

2. Forecast step:

Input: \mathbf{x}^a

Output: \mathbf{x}^f

3. Ensemble selection:

3.1. Anomalies generation

- Select an ensemble \mathbf{X}
- Compute the anomalies $\mathbf{X}' = \mathbf{X} - \bar{\mathbf{x}}$ where $\bar{\mathbf{x}}$ is the mean of \mathbf{X}

3.2. \mathbf{X}^f generation

Inputs: \mathbf{X}' and \mathbf{x}^f

Output: $\mathbf{X}^f = \mathbf{X}' + \mathbf{x}^f$

4. Goto 1

Algorithm 2 Generic ensemble design Algorithm

1. Inputs: a dictionary $\mathbf{D} = [\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_L]$ of model outputs, the desired ensemble size N (with $L \gg N$ and at least $L \geq N$), and the forecast \mathbf{x}^f is iterated through the dictionary to apply the selection.
 2. Sort the elements based on the metric ordering criteria:
 $\mathbf{d}_{j_1}, \mathbf{d}_{j_2}, \dots, \mathbf{d}_{j_N}, \dots, \mathbf{d}_{j_L}$
 3. Form the ensemble of the first N members $\mathbf{X} = [\mathbf{d}_{j_1}, \mathbf{d}_{j_2}, \dots, \mathbf{d}_{j_N}]$ and use it to update the forecast with the incoming observations.
-

Table 1: Specific algorithms of the selection methods

Selection method	Algorithm
L2-norm	<p>2. For i from 1 to L compute $\gamma_i = \ \mathbf{x}^f - \mathbf{d}_i\ _2$</p> <p>3. Sort the γ_i in ascending order: $\gamma_{i_1}, \gamma_{i_2}, \dots, \gamma_{i_L}$ with $\gamma_{i_1} \leq \gamma_{i_2} \leq \dots \leq \gamma_{i_L}$ and assign $j_1 = i_1, j_2 = i_2, \dots, j_L = i_L$</p>
OMP	<p>2.1 Initialization: set $\mathbf{y}_0 = \mathbf{0}$, index set $\Delta_0 = \emptyset$ and residual $\mathbf{r}_0 = \mathbf{x}^f$</p> <p>2.2 For t from 1 to N,</p> <ul style="list-style-type: none"> • Find the index of the dictionary element having the highest inner product with the residual: set δ_t to one of the indexes j for which the maximum is reached, i.e. $\langle \mathbf{r}_{t-1}, \mathbf{d}_{\delta_t} \rangle = \max_{j=1, \dots, L} \langle \mathbf{r}_{t-1}, \mathbf{d}_j \rangle$ • Augment the index set: $\Delta_t = \Delta_{t-1} \cup \{\delta_t\}$ • Solve the least-square problem $\min_y \ \mathbf{x}^f - \mathbf{D}_{\Delta_t} \mathbf{y}\ _2$ and then choose $\mathbf{y}_t \in \arg \min_y \ \mathbf{x}^f - \mathbf{D}_{\Delta_t} \mathbf{y}\ _2$ • Calculate new residual $\mathbf{r}_t = \mathbf{x}^f - \mathbf{D}_{\Delta_t} \mathbf{y}_t$ <p>End for</p> <p>3. Assign $j_1 = \delta_1, j_2 = \delta_2, \dots, j_N = \delta_N$</p>

Table 2: Region wise SSH statistics for CMEMS-L4 interpolated product, *Fexp*, EAKF, EnOI, AEnOI-L2, and AEnOI-OMP. Statistics are shown for four different regions, Gulf of Aden (GoA; 30°E-50°E and 10°N-14°N), Southern Red Sea (SRS; 30°E-50°E and 14°N-19°N), Central Red Sea (CRS; 30°E-50°E and 19°N-23°N) and Northern Red Sea (NRS; 30°E-50°E and 23°N-28°N). Units for mean, STD and RMSD are in cm. The assimilation experiment yielding best results for a region is highlighted with bold fonts.

	GoA				SRS				CRS				NRS			
	Mean	STD	RMSE	Corr	Mean	STD	RMSE	Corr	Mean	STD	RMSE	Corr	Mean	STD	RMSE	Corr
Observation	12	5			12	11			5	13			-5	13		
CMEMS-L4	7	5	5	0.94	7	11	5	0.98	1	12	5	0.98	-10	13	6	0.98
<i>Fexp</i>	5	5	10	0.40	6	8	9	0.84	-1	10	9	0.86	-13	10	11	0.82
EAKF	13	6	5	0.57	14	11	6	0.88	8	12	7	0.86	-4	11	7	0.85
EnOI	13	8	7	0.56	12	11	5	0.91	6	13	6	0.90	-6	12	6	0.89
AEnOI-L2	11	7	6	0.61	12	11	5	0.92	7	13	6	0.89	-4	12	6	0.88
AEnOI-OMP	10	8	7	0.54	12	11	5	0.88	7	13	7	0.88	-4	11	7	0.85

List of Figures

1	Workflow of the dictionary-based AEnOI schemes as implemented with DART. The forecast state \mathbf{x}^f is used to select an ensemble $\mathbf{X}^f = [\mathbf{x}_1^f, \dots, \mathbf{x}_N^f]$ from an available dictionary. This ensemble is then centered around \mathbf{x}^f before it is updated by the upcoming observation to obtain the analysis state \mathbf{x}^a , which is then integrated by the model to compute the next forecast.	37
2	Illustration of an ensemble construction based on L2. Compute the L2-distances ($dist_1, dist_2, \dots, dist_L$) between the forecast x^f and the dictionary members (d_1, d_2, \dots, d_L) then select the first N members ($d_{j_1}, d_{j_2}, \dots, d_{j_N}$) with the smallest distances to the forecast member to generate the ensemble X	38
3	Illustration of an ensemble construction based on OMP. Compute the inner products (ip_1, ip_2, \dots, ip_L) between the forecast x^f and the dictionary members (d_1, d_2, \dots, d_L) and keep the member having the highest ip value. Solve the least-square problem between the forecast and that member, and then compute the residual r_1 . Compute the inner products between the residual r_1 and the remaining dictionary members and keep the member having the highest ip value. Solve the least-square problem between the forecast and the set containing that member and all the previous selected members. Compute the residual r_2 . Repeat the process with the successive residuals until N members are selected.	39
4	Mean analysis RMSE (top) and forecast ensemble standard deviation (bottom) for EnKF, EnOI, AEnOI-L2 and AEnOI-OMP schemes using the L-63 model. The assimilation experiments were performed using 100 members and assimilating all data every 4 model steps.	40
5	Sensitivity of the ensemble schemes to the ensemble size for a given assimilation window using L-63 model.	41
6	Sensitivity of the ensemble schemes to the assimilation period for a given ensemble size using L-63 model.	42

7	Time-averaged RMSE as a function of the localization radius (x axis) and inflation factor (y axis) using L-96 model. All filters were implemented with 40 members, and observations were assimilated every 4 model time steps. The minimum RMSEs are indicated by asterisks, and their associated values are given in the title.	43
8	Sensitivity of the ensemble schemes to the ensemble size for a given assimilation window using L-96 model.	44
9	Sensitivity of the ensemble schemes to the length of the assimilation period for a given ensemble size using L-96 model.	45
10	Spatial maps of SST STD in °C (b-g), Bias (h-m), RMSE (n-s) and correlations (t-y) for OSTIA (b), <i>Fexp</i> (c), EAKF (d), EnOI (e), AEnOI-L2 (f), and AEnOI-OMP (g). All the statistics are with respect to satellite level-3 SST observations. Panel “a” shows STD in the satellite level-3 SST. Negative values of bias indicate model cold biases and vice versa.	46
11	Time series of root-mean-square-error (RMSE) for daily averaged forecasts of (a) SST (b) SSH from <i>Fexp</i> (red), EAKF (maroon), EnOI (green), AEnOI-L2 (blue), AEnOI-OMP (pink), and level-4 gridded products (OSTIA for SST and CMEMS-L4 for SSH; black). RMSE is computed by collocating the daily averaged model forecasts onto satellite along-track level-3 SST and SSH observations. 10-day smoothing is applied to better emphasize the differences between the assimilation results. Units are in “°C” and “cm” for SST and SSH, respectively. Panels (c) and (d) are similar to (a) and (b) except that the RMSEs are computed for 3-day spaced analyses (snapshots after assimilation) without smoothing.	47
12	Depth-Time evolution of temperature (°C; a, c, e, g, and i) and salinity (psu; b, d, f, h, and j) and depth of 22 °C isotherm (meters; k) at (38°E, 22°N) as resulted from <i>Fexp</i> , EAKF, EnOI, AEnOI-L2, and AEnOI-OMP.	48
13	Horizontal and vertical distributions of ensemble spread (a - d) of SSH, (e - h) SST, and (i - l) temperature on 1-May-2011 as they result from EAKF, EnOI, AEnOI-L2, and AEnOI-OMP.	49

14	The histograms (prior) in experiments using EnOI (1 st column), AEnOI-L2 (2 nd column) and AEnOI-OMP (3 rd column) assimilation experiments at three selected locations (indicated in Figure 15) in the northern, central and southern basins of the Red Sea, as the 1 st , 2 nd and 3 rd rows, respectively.	50
15	Sampled correlations for SST as computed from the assimilation experiments using EnOI (1 st column), AEnOI-L2 (2 nd column), and AEnOI-OMP (3 rd column) schemes at three selected locations in the northern (1 st row), central (2 nd row), and southern (3 rd row) basins of the Red Sea, respectively, before applying localization. The black dot in each panel indicates the selected location.	51

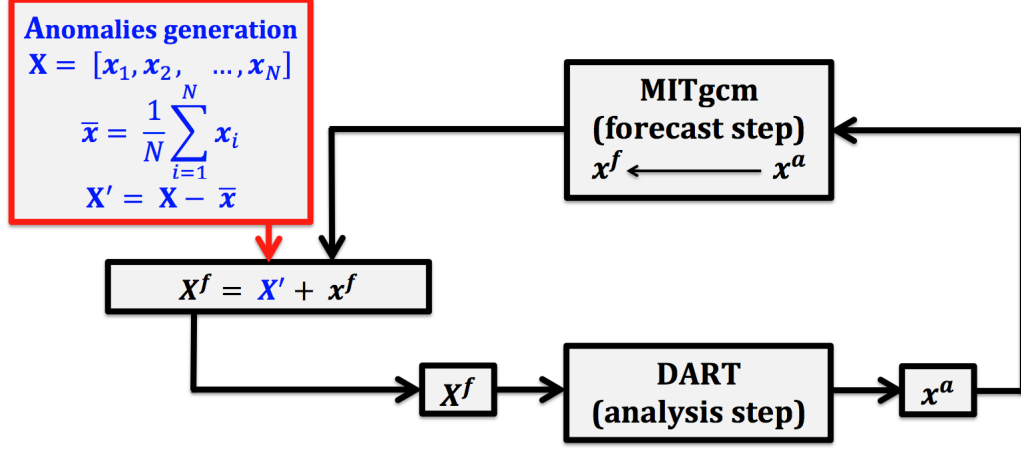


Figure 1: Workflow of the dictionary-based AEnOI schemes as implemented with DART. The forecast state \mathbf{x}^f is used to select an ensemble $\mathbf{X}^f = [\mathbf{x}_1^f, \dots, \mathbf{x}_N^f]$ from an available dictionary. This ensemble is then centered around \mathbf{x}^f before it is updated by the upcoming observation to obtain the analysis state \mathbf{x}^a , which is then integrated by the model to compute the next forecast.

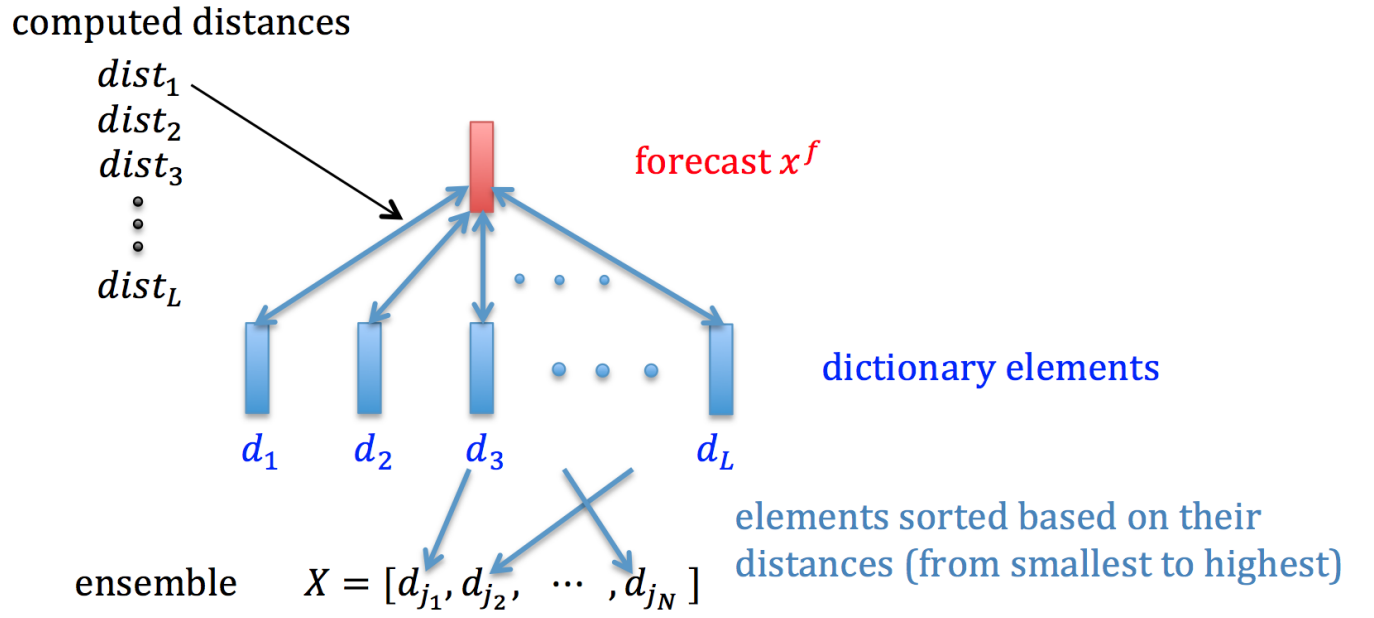


Figure 2: Illustration of an ensemble construction based on L2. Compute the L2-distances ($dist_1, dist_2, \dots, dist_L$) between the forecast x^f and the dictionary members (d_1, d_2, \dots, d_L) then select the first N members ($d_{j_1}, d_{j_2}, \dots, d_{j_N}$) with the smallest distances to the forecast member to generate the ensemble X .

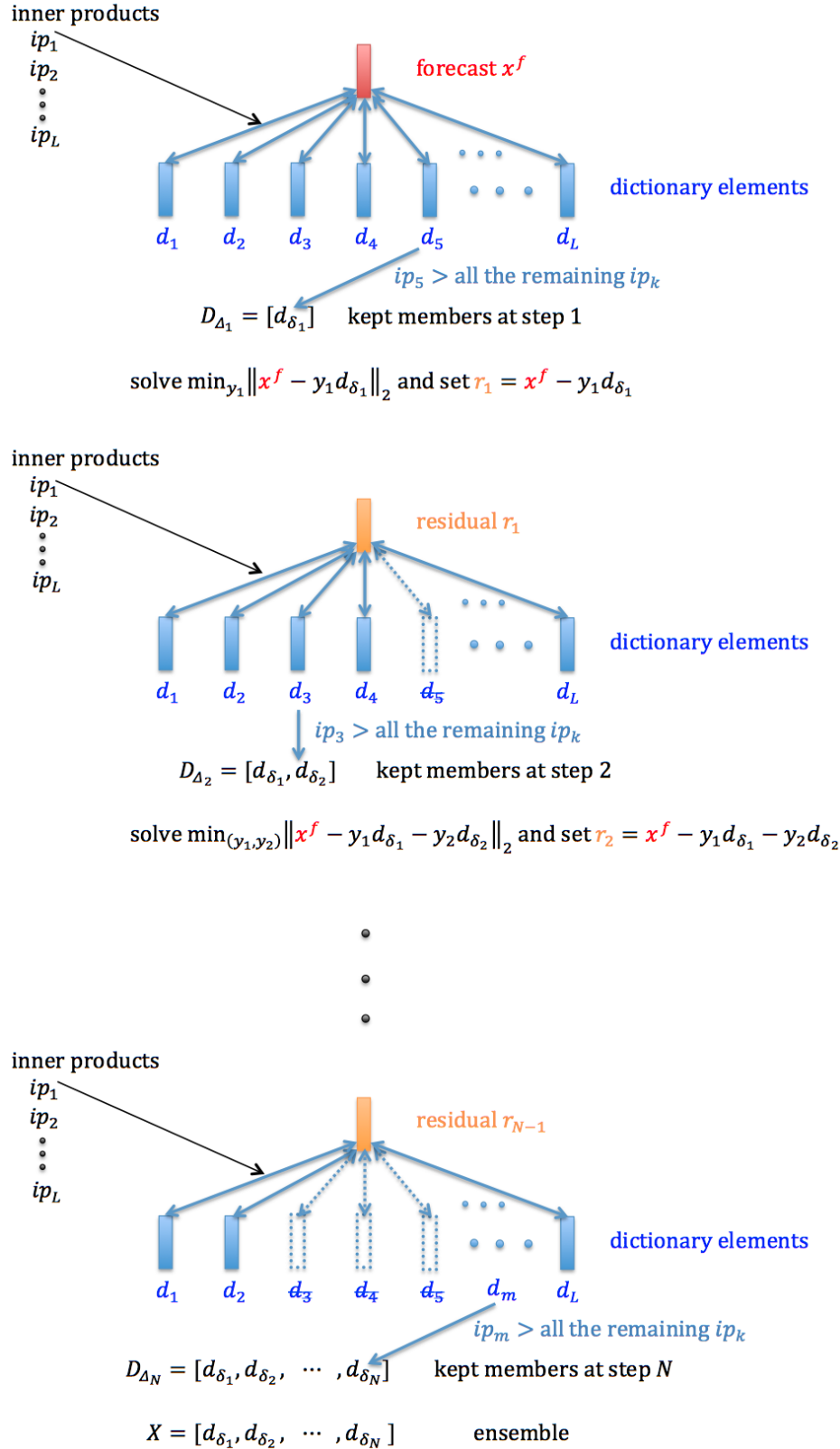


Figure 3: Illustration of an ensemble construction based on OMP. Compute the inner products (ip_1, ip_2, \dots, ip_L) between the forecast x^f and the dictionary members (d_1, d_2, \dots, d_L) and keep the member having the highest ip value. Solve the least-square problem between the forecast and that member, and then compute the residual r_1 . Compute the inner products between the residual r_1 and the remaining dictionary members and keep the member having the highest ip value. Solve the least-square problem between the forecast and the set containing that member and all the previous selected members. Compute the residual r_2 . Repeat the process with the successive residuals until N members are selected.

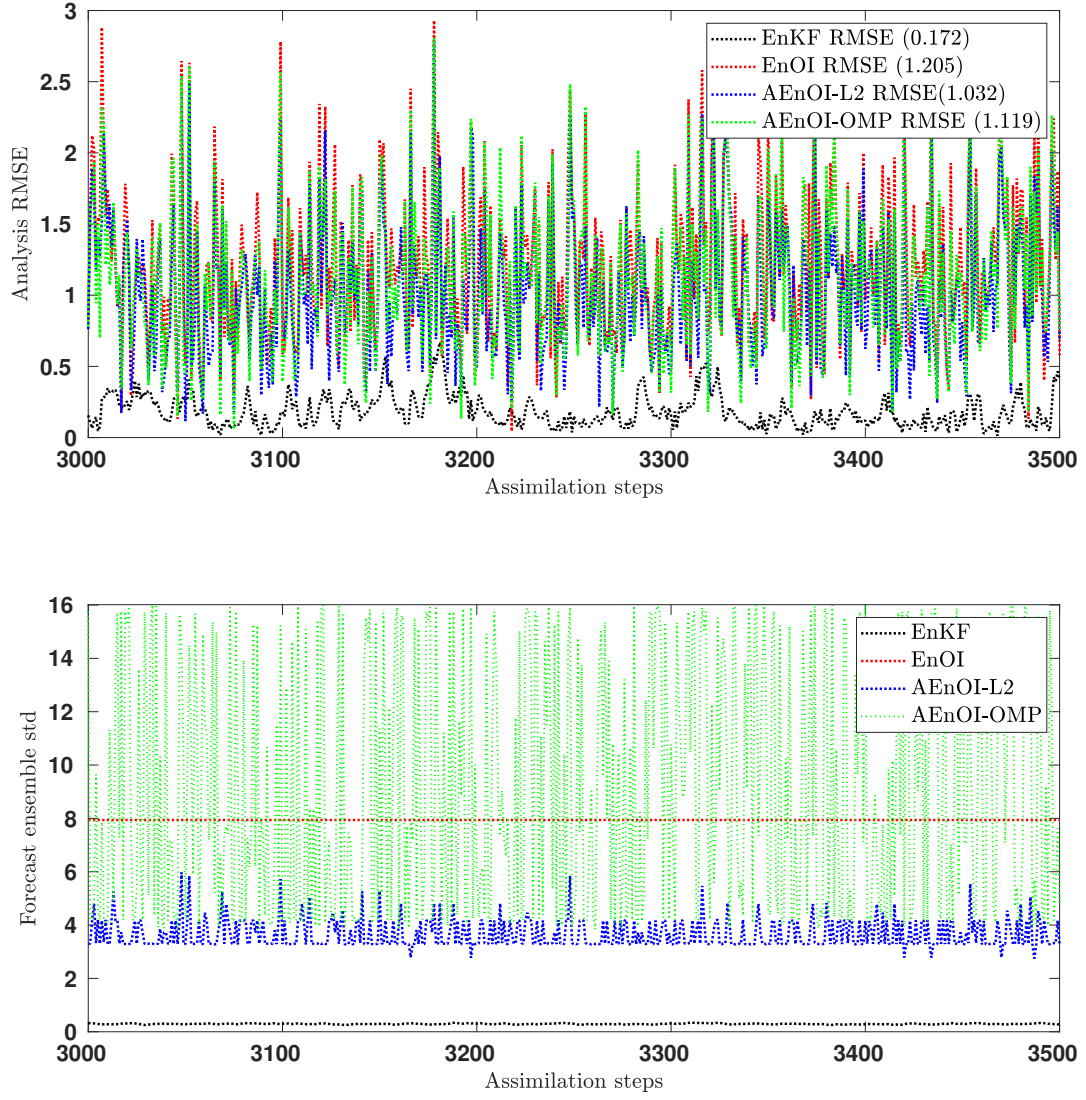


Figure 4: Mean analysis RMSE (top) and forecast ensemble standard deviation (bottom) for EnKF, EnOI, AEnOI-L2 and AEnOI-OMP schemes using the L-63 model. The assimilation experiments were performed using 100 members and assimilating all data every 4 model steps.

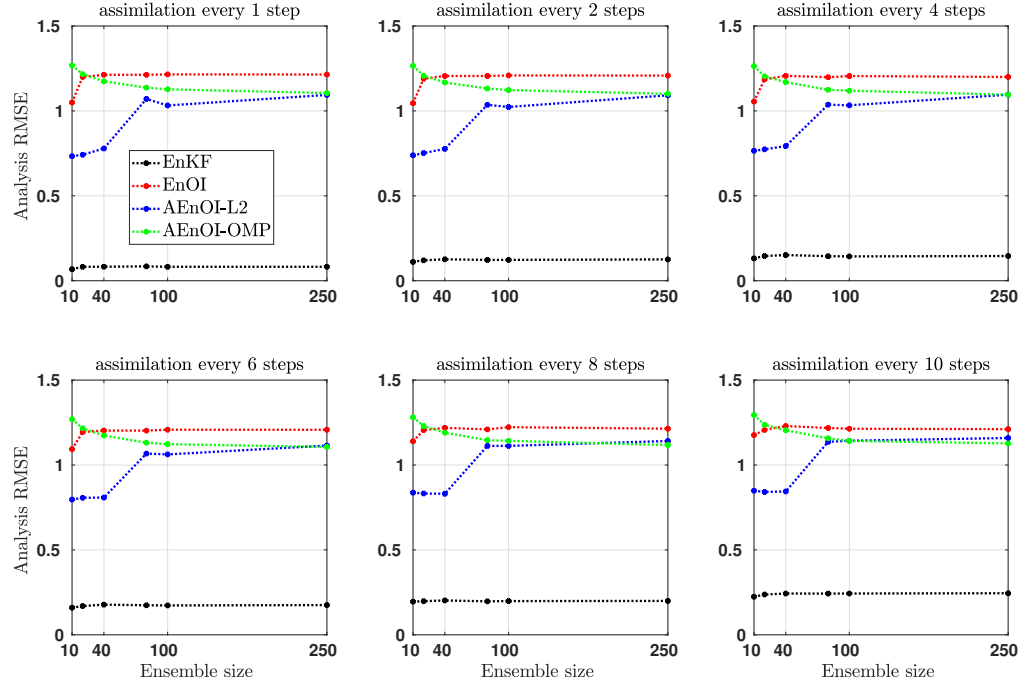


Figure 5: Sensitivity of the ensemble schemes to the ensemble size for a given assimilation window using L-63 model.

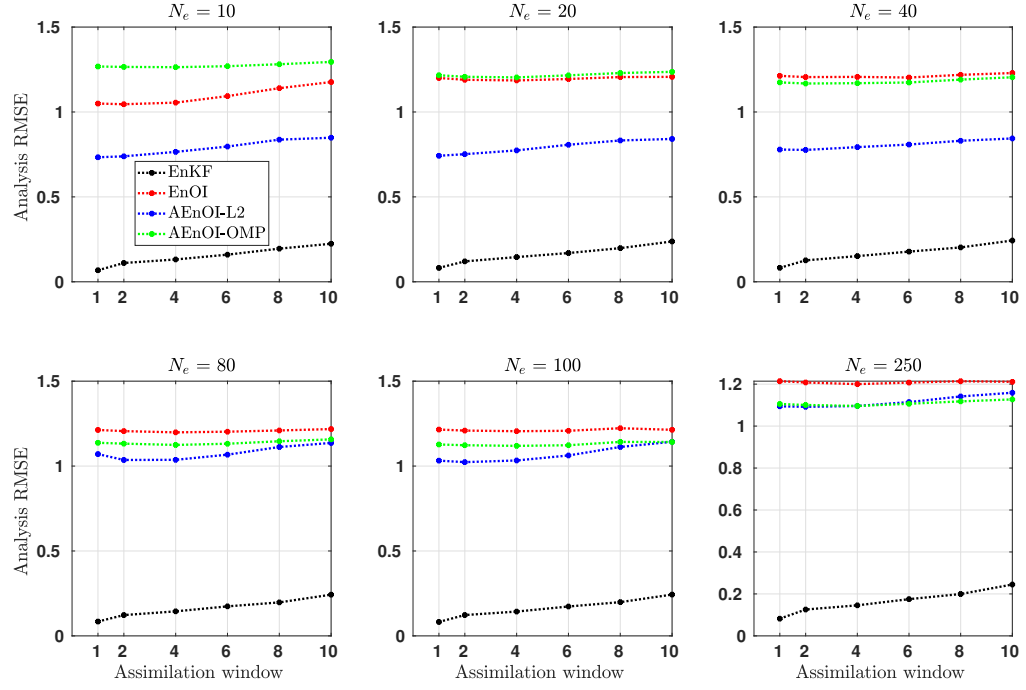


Figure 6: Sensitivity of the ensemble schemes to the assimilation period for a given ensemble size using L-63 model.

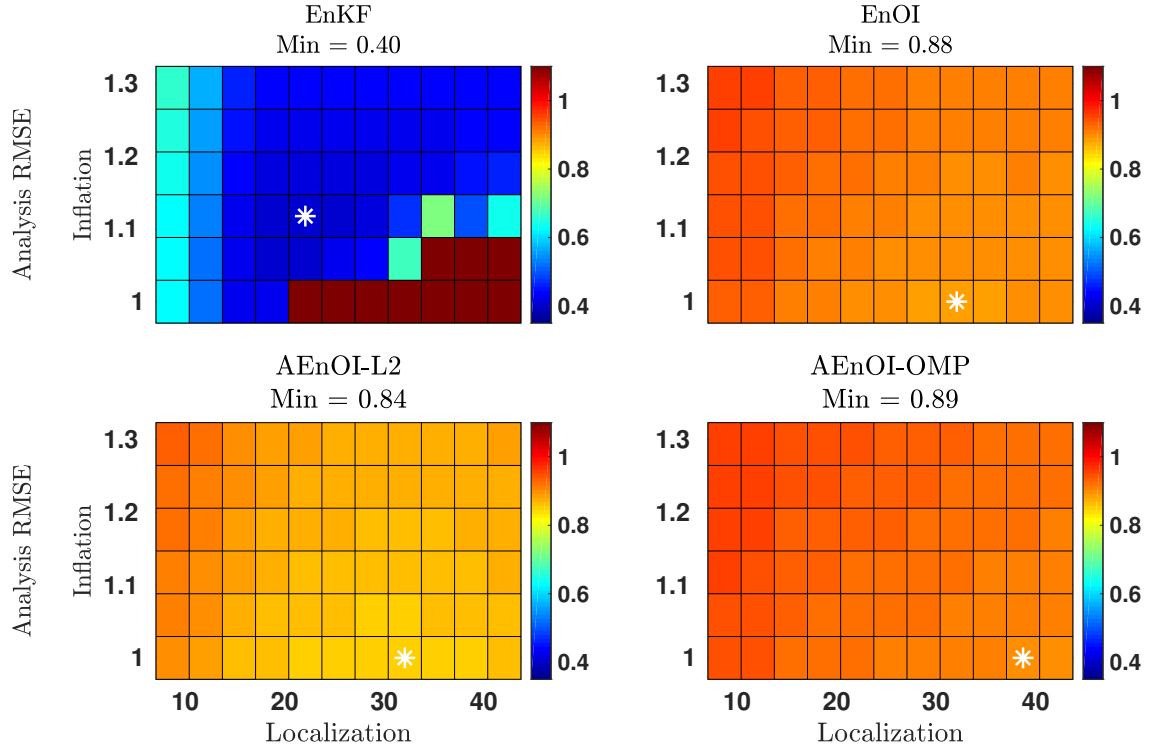


Figure 7: Time-averaged RMSE as a function of the localization radius (x axis) and inflation factor (y axis) using L-96 model. All filters were implemented with 40 members, and observations were assimilated every 4 model time steps. The minimum RMSEs are indicated by asterisks, and their associated values are given in the title.

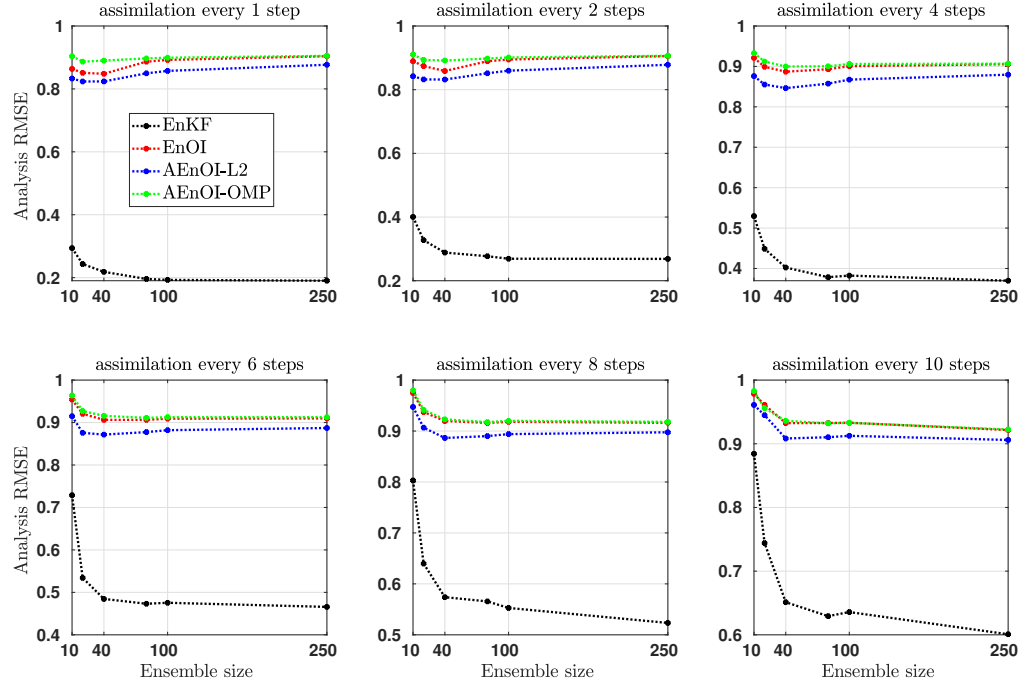


Figure 8: Sensitivity of the ensemble schemes to the ensemble size for a given assimilation window using L-96 model.

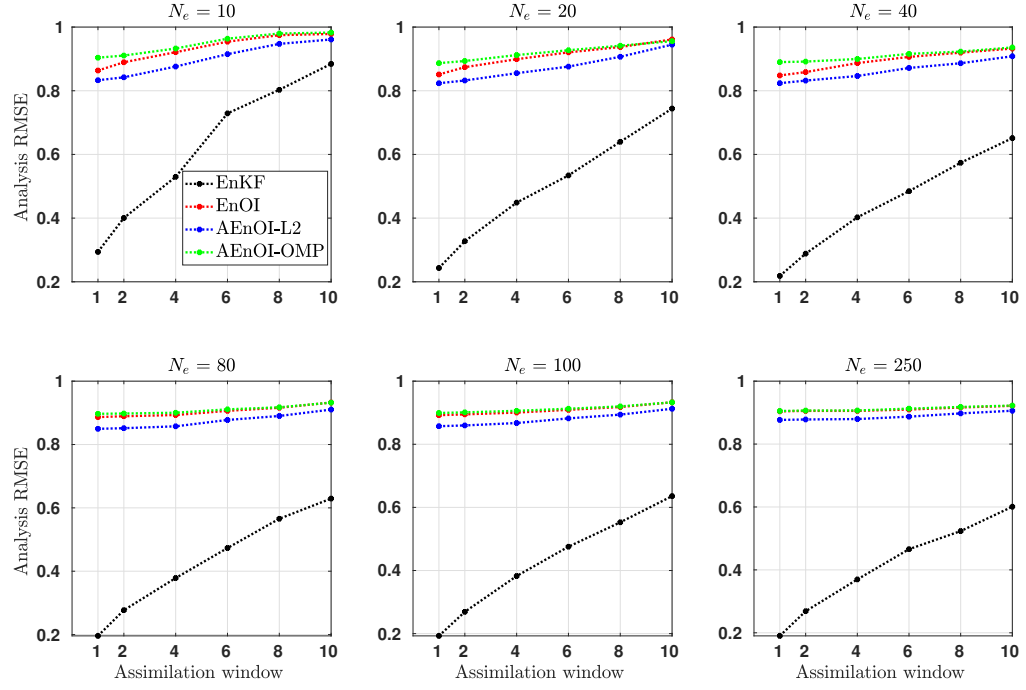


Figure 9: Sensitivity of the ensemble schemes to the length of the assimilation period for a given ensemble size using L-96 model.

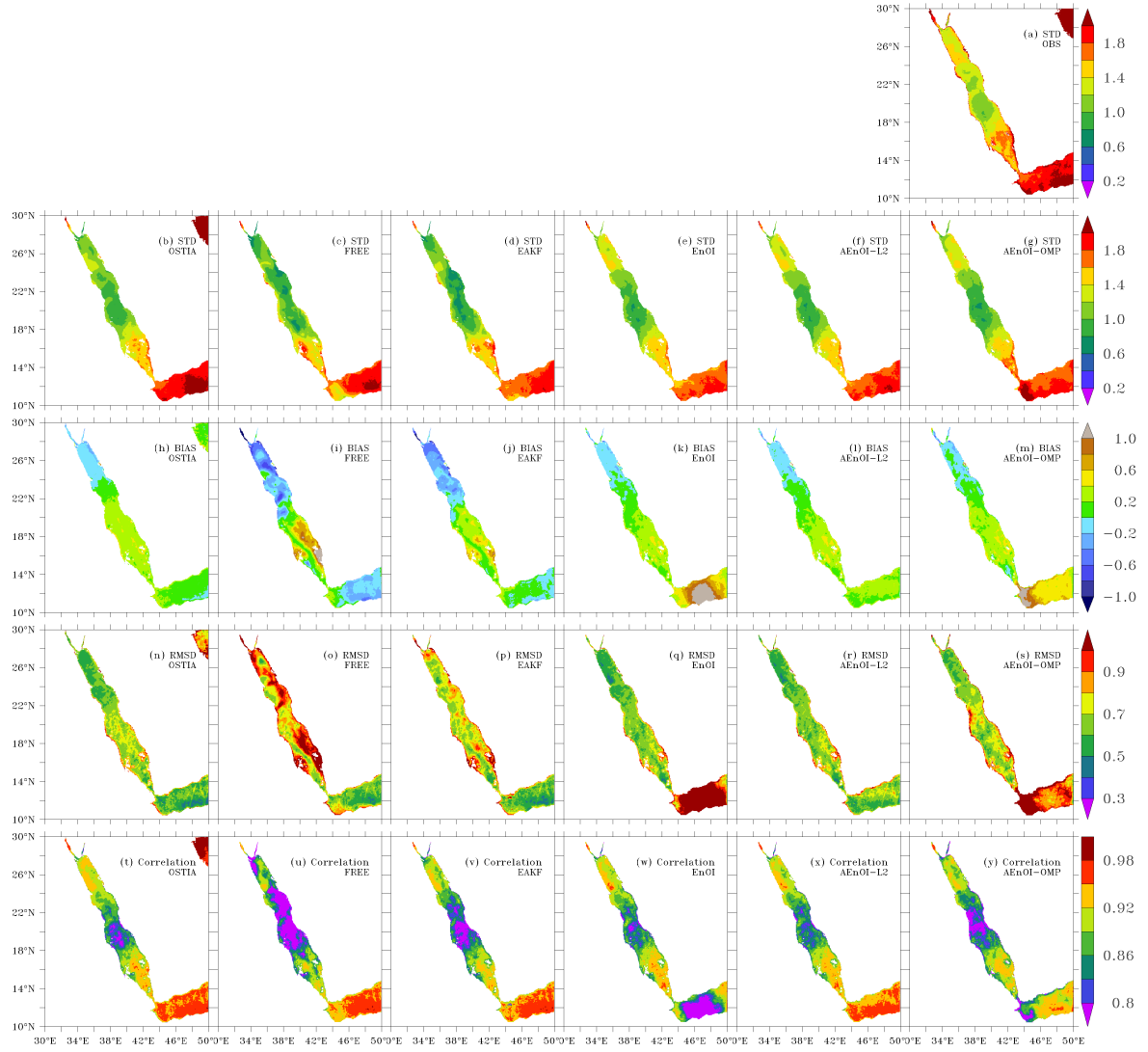


Figure 10: Spatial maps of SST STD in $^{\circ}\text{C}$ (b-g), Bias (h-m), RMSE (n-s) and correlations (t-y) for OSTIA (b), *Free* (c), EAKF (d), EnOI (e), AEnOI-L2 (f), and AEnOI-OMP (g). All the statistics are with respect to satellite level-3 SST observations. Panel “a” shows STD in the satellite level-3 SST. Negative values of bias indicate model cold biases and vice versa.

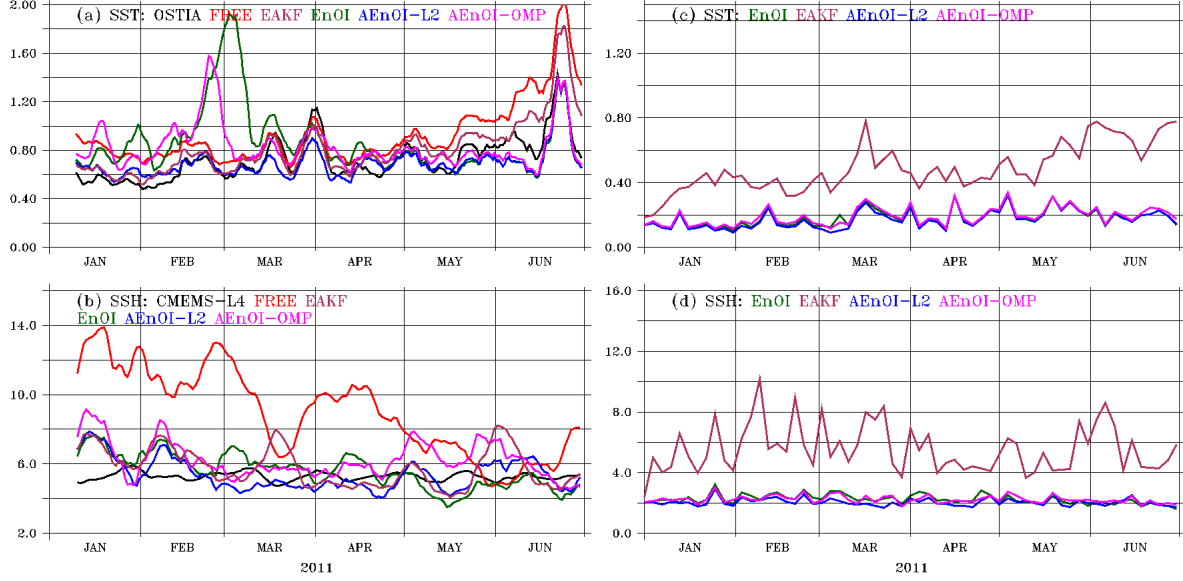


Figure 11: Time series of root-mean-square-error (RMSE) for daily averaged forecasts of (a) SST (b) SSH from *Exp* (red), EAKF (maroon), EnOI (green), AEnOI-L2 (blue), AEnOI-OMP (pink), and level-4 gridded products (OSTIA for SST and CMEMS-L4 for SSH; black). RMSE is computed by collocating the daily averaged model forecasts onto satellite along-track level-3 SST and SSH observations. 10-day smoothing is applied to better emphasize the differences between the assimilation results. Units are in “°C” and “cm” for SST and SSH, respectively. Panels (c) and (d) are similar to (a) and (b) except that the RMSEs are computed for 3-day spaced analyses (snapshots after assimilation) without smoothing.

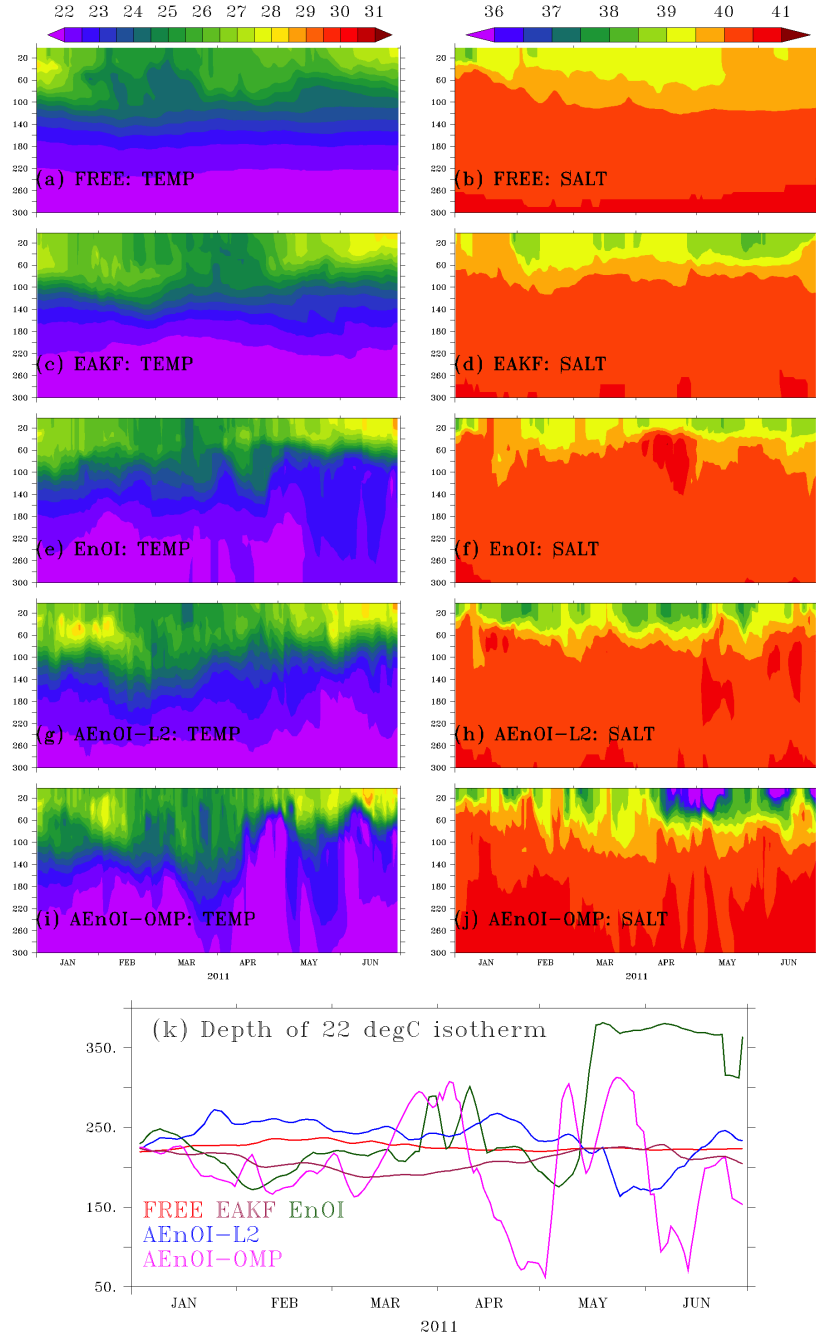


Figure 12: Depth-Time evolution of temperature ($^{\circ}\text{C}$; a, c, e, g, and i) and salinity (psu; b, d, f, h, and j) and depth of 22 $^{\circ}\text{C}$ isotherm (meters; k) at (38°E , 22°N) as resulted from *Free*, EAKF, EnOI, AEnOI-L2, and AEnOI-OMP.

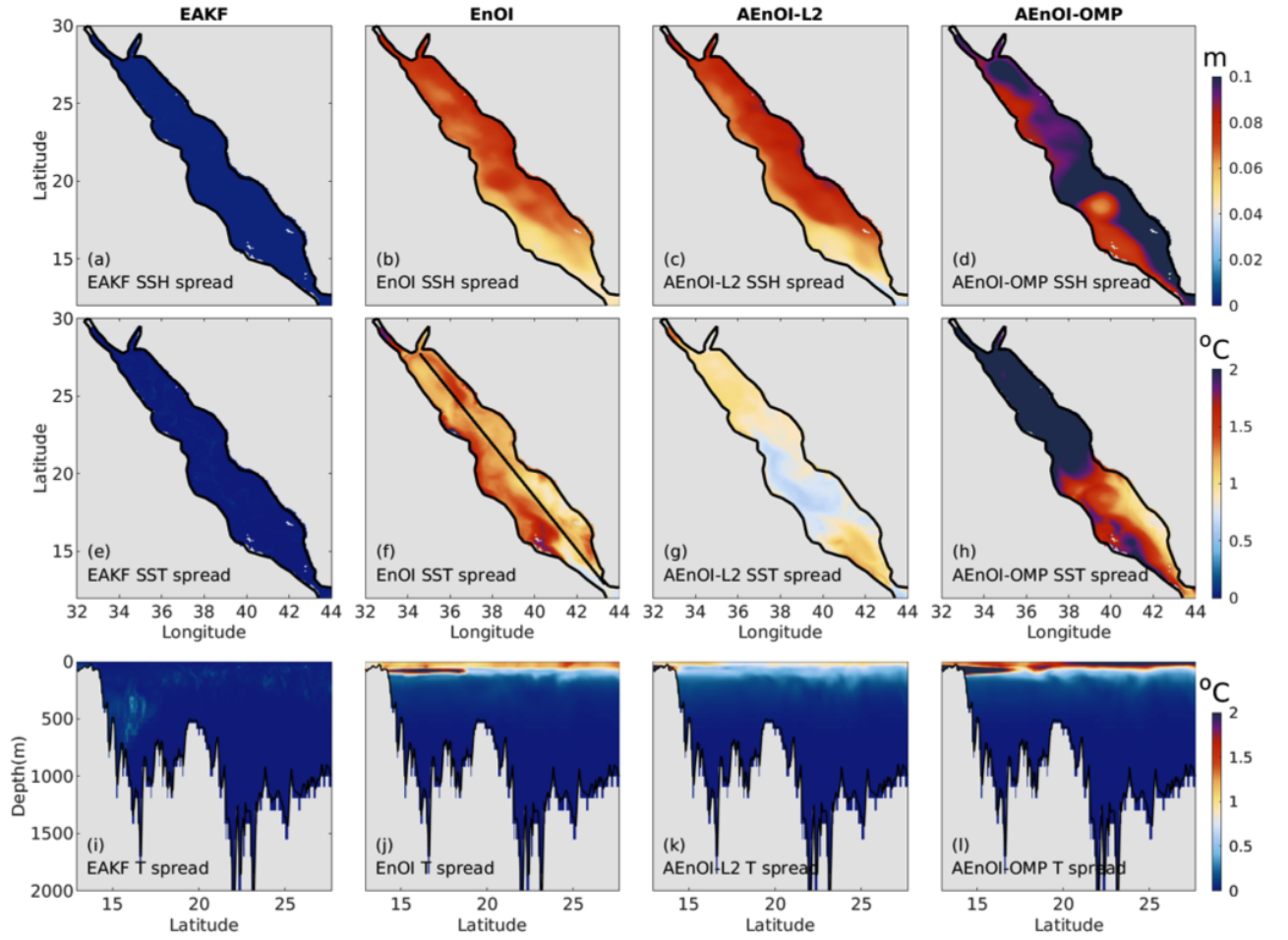


Figure 13: Horizontal and vertical distributions of ensemble spread (a - d) of SSH, (e - h) SST, and (i - l) temperature on 1-May-2011 as they result from EAKF, EnOI, AEnOI-L2, and AEnOI-OMP.

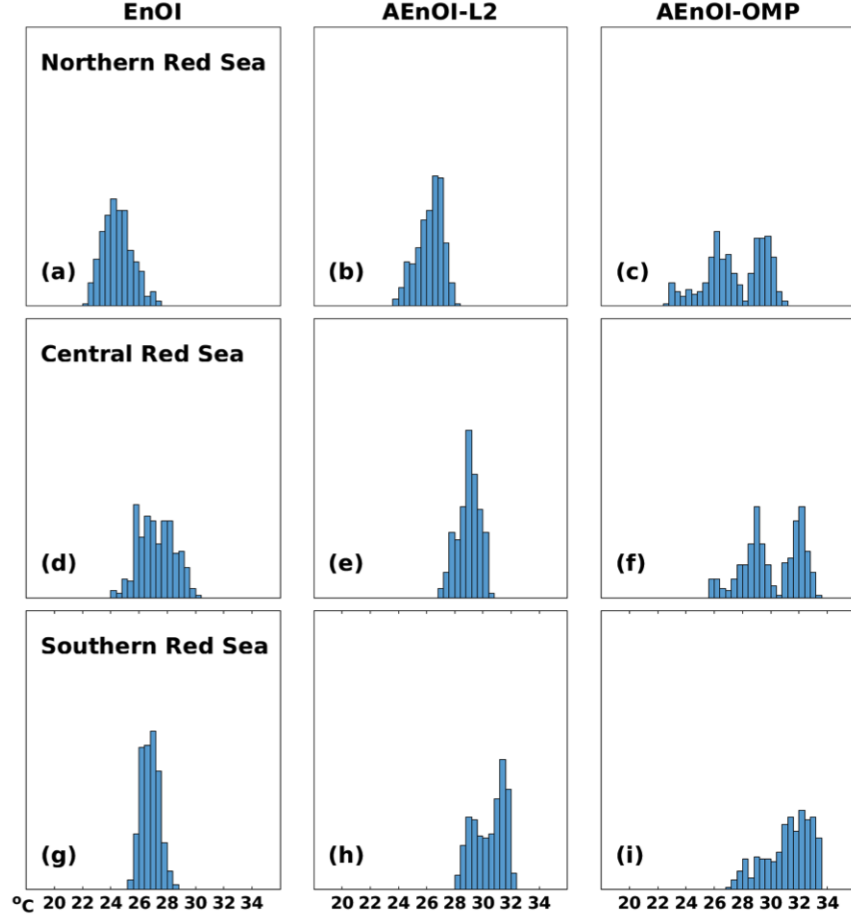


Figure 14: The histograms (prior) in experiments using EnOI (1st column), AEnOI-L2 (2nd column) and AEnOI-OMP (3rd column) assimilation experiments at three selected locations (indicated in Figure 15) in the northern, central and southern basins of the Red Sea, as the 1st, 2nd and 3rd rows, respectively.

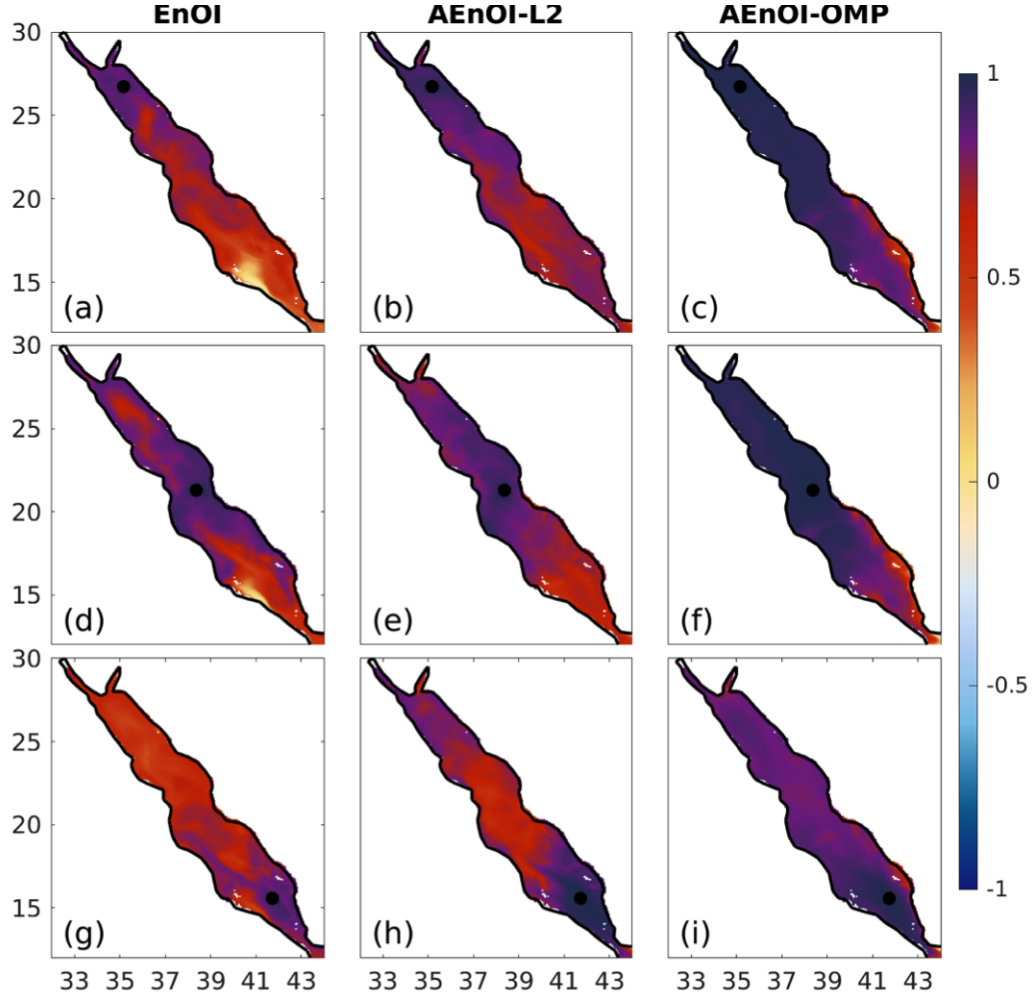


Figure 15: Sampled correlations for SST as computed from the assimilation experiments using EnOI (1st column), AEnOI-L2 (2nd column), and AEnOI-OMP (3rd column) schemes at three selected locations in the northern (1st row), central (2nd row), and southern (3rd row) basins of the Red Sea, respectively, before applying localization. The black dot in each panel indicates the selected location.