# Generalized data assimilation method and its comparison with an Ensemble Optimal Interpolation scheme in conjunction with the numerical ocean model using altimetry data

**Konstantin Belyaev[1,2], Andrey Kuleshov[2], Ilya Smirnov[3] and Clemente A. S. Tanajura[4]**

[1]Shirshov Institute of Oceanology of Russian Academy of Sciences, Moscow, 117997, Russia
[2]Keldysh Institute of Applied Mathematics of Russian Academy of Sciences, Moscow, 125047, Russia
[3]Lomonosov Moscow State University, Faculty of Computational Mathematics and Cybernetics, Moscow, 119991, Russia
[4]Federal University of Bahia, Physics Institute and Center for Research in Geophysics and Geology, Salvador BA, 40170-280, Brazil

**Correspondence:** Andrey Kuleshov (andrew_kuleshov@mail.ru)

**Abstract.** A recently developed original data assimilation scheme is presented and tested. The scheme is based on the application of the theory of diffusion random processes. It is applied here in conjunction with the Hybrid-Coordinate Ocean Model (HYCOM) to assimilate altimetry data from the Archiving, Validating and Interpolating Satellite Oceanography Data (AVISO) in the Atlantic. Several numerical experiments were carried out and their results were analyzed. It is shown that the method can assimilate data and provide analyses closer to observations in a sense of standard metrics. The quality of the data assimilation can be estimated due to the results which becomes close to the observations then the model run without assimilation. This "closeness" is numerically assessed by different metrics in particular, by the variance of the error between data and model values at the same points. This method allows calculating the confidence range of the analyses by estimating their errors. The presented method is compared with the Ensemble Optimal Interpolation scheme (EnOI) and it is shown that it has several advantages, in particular, it provides a better forecast and requires less computational cost.

## 1 Introduction

Nowadays, data assimilation (DA) theory, as a part of mathematical and numerical research, is a scientific area of a great practical importance. The ideas and methods of DA are used in ocean modelling, weather forecasting, operational oceanography and many other fields of science. Its main goal is to combine the numerical model results with independently observed data in a mathematically optimal way to represent the physical system state. Since the two extreme approaches for this representation are obviously poor - namely, they ignore observations altogether or replace only the observed model variables with observational data without changing the others – the optimal way of producing the optimal representation of reality is not trivial.

It is important to point out that the terms "optimal", "combine", etc. do not make much sense unless they are defined within a strong theoretical framework. Therefore, in order to give meaning to these formulations, it is necessary to consider the data assimilation problem as a part of one or several mathematical and/or physical theories, such as the optimal control theory, mathematical statistics, numerical analysis, etc.

The studies dealing with the representation of physical systems' behavior were carried out and published in the scientific literature for more than 50 years, since the beginning of the 1960s. A good review of the main achievements in this research domain over the last century is presented in Ghil et al. (1991). Since the beginning of the 2000s, the main progress is related to the development of computer facilities, explosion of the observational data network, parallel computations, and other technical novelties. This advance led to the progress in new mathematical methods and algorithms, construction and development of numerical models with a very high resolution, data exchange all over the world, etc. At the present time, the DA techniques, algorithms, and methods have become an essential part of operational oceanography in the ocean shelf and

coastal areas, especially in the oil and gas mining zones, as well as in the zones of pipeline transportation. Several national and international scientific projects are aimed to search for the optimal solution of the DA techniques in conjunction with regionally configured numerical models. In particular, one can name the Brazilian REMO project (Lima et al., 2016), the Australian Blue Link project (Schiller et al., 2011), the American HYCOM&NCODA project (Cummings et al., 2013) and others.

However, the development and application of new and more advanced DA schemes and methods remains a topical and very important theoretical and practical problem. There is a necessity to have a powerful and, at the same time, relatively portable and economy DA scheme which would be applicable to various numerical ocean and coupled ocean-atmosphere models and would provide a satisfactory and reliable forecast of the ocean characteristics for short- and medum-term periods.

Many of the used DA methods and approaches may be divided into two large groups. The first group is based on the variational or functional schemes; its modern version is known in literature as the 4D_VAR approach. Its main idea consist in finding such a model trajectory within a given time window that this trajectory passes nearby the observations as close as possible represented by the minimum of a known functional. In fact, this approach is reduced to solving the inverse problem to find the initial and boundary conditions of the model and obtain the model trajectory that starts from the sought fields close to observations with respect to a given metric of closeness. Normally, one should minimize some functional as a function of the sought trajectory in a known metric, which is given a priori. Several recent investigations in this direction can be found in (e.g. Marchuk et al., 2012; Talagrand et al., 1987; Agoshkov et al., 2010).

The second large group is known as a "dynamical-stochastic" DA approach. This approach uses the theory of statistical estimation or filtration to find out the best estimator to find out the best estimator to minimize the variance relatively the observations of this new constructed field among all others theoretically possible fields. Usually, the observed variables are originally supposed to be a sum of a true signal., which the model is supposed to represent, and stochastic noise with known probabilistic properties. Started from the so-called "objective interpolation method" (see for instance Pendruff et al., 2002), its modern development is known in the literature as the Ensemble Kalman filter method (EnKF). As the examples one may refer to the papers of Evensen (2009) and Xie et al. (2010).

There are also, relatively small numbers of papers where the considered DA methods distinguish from that two groups. For instance, in Van Leeuwen (2015) the corresponding scheme based on the partial method is applied. Also, in Van Leeuwen, (2011) the well-known Bayes technique is used to find the best estimation of a posteriori probability after one time-step model simulation. However, the most actually used DA techniques relate to the schemes indicated before.

Recently hybrid DA schemes combining both approaches have appeared. In those schemes both ideas, namely, a functional and dynamical-stochastic are exploited together. In particular, the minimum of functional is sought in the metric which represents the variance of the stochastic variable or variables that are determined from the observations. Some examples may be found (e.g. Lorenc et al., 2015; Tanajura et al., 2013). We may also refer to Tanajura et al. (2009), where the hybrid DA method based on theory of diffusion stochastic processes was proposed.

This study uses the ideas presented in Tanajura et al. (2009), however dealing with the application of the novel DA method created in Belyaev et al. (2018), Belyaev et al. (2019). Hereafter, this method will be referred to as the GKF, or the Generalized Kalman Filter. Unlike the standard EnKF, which considers only the model output, called background, and observed data at the assimilation moment, the GKF considers additionally the temporal tendency from both the model and observations. The explicit form of the obtained gain matrix (analogous of the Kalman gain matrix in the standard scheme) contains the time-derivative of the average of the model ensemble fields in two sequential time-moments. As a consequence, the gain matrix proves to be zero if these tendencies coincide with each other and, thus, no assimilation occurs. The GKF becomes better in a sense of numerical forecast of the ocean state than its counterpart, EnKF in a case the assimilation occurs since it takes into account not only the instantaneous data but also the time-derivative. The mathematical derivation of the

GKF is not trivial., and the paper Belyaev et al. (2018) specially is dedicated to its formulation and mathematical proofs, but its physical sense is rather transparent. The method is based on the undeniable physical axiom, namely, on the path-of-least-resistance principle. According to this principle, all transitions from one physical state to other pass with the minimum energy change, which may be set by the specific Lagrange functional. Since actually, the assimilation process is the transform from the model ocean state without specific data to the model state including the new coming information, namely the information from data, the corresponding transition must take place and this transition is realized with respect to this principle. Paper Belyaev et al. (2018) contains all the mathematical aspects of these ideas, but this study is focused on its application, feasible realizations and analysis of the results.

Here, this method is used in conjunction with the ocean model HYCOM, presented in Bleck (2002), Bleck et al. (1981). In the current work, this method is applied to the dynamical simulation in the Atlantic and is compared with the standard Ensemble Optimal Interpolation (EnOI) method (Evensen, 2009), a simplified version of the EnKF, as an alternative data assimilation scheme. Twin experiments were performed for the same initial conditions with the assimilation of the equivalent data. The AVISO data (www.aviso.org) was chosen as the input data for the assimilation.

Concisely, it is possible to point out how two or more data assimilation methods may be compared. There are several criteria, including computational consumptions, feasibility of their realization, reliability of their implementations etc. However, the most obvious criterion of comparison may be determined as the minimum variance of the forecast error. Indeed, the basic idea of using any DA technique is to minimize the model forecast error if the model starts from the corrected field after assimilation. If one data assimilation method does it better then another, it is obviously preferable. Such a comparison criterion was earlier used in Belyaev et al. (2012), where the standard method of objective interpolation, OI was compared with the EnKF and method which is described in Tanajura and Belyaev (2009). Similarly, the OI method was compared with the EnOI accordingly to this criterion in Kaurkin et al. (2018).

The main goals of our study are as follows: (i) to present the feasibility and applicability of the referred GKF method, including its parallelization and computational effectiveness; (ii) to compare this method with the alternative EnOI scheme and prove that the GKF has many advantages, including lower computational consumptions and better forecast properties; (iii) to analyze the results of modeling and show that the GKF captures the basic structure of synoptic variability in the Atlantic; (iv) to estimate the analysis error appeared in the model assimilation and its dynamics in time.

The structure of this work is as follows: Section 1 is the introduction, Section 2 is the description of the data assimilation method and the algorithm of its realization, Section 3 is the analysis of the results and the comparison with the EnOI method and with the control, i.e. the model simulation with the same initial conditions and forcing but without data assimilation. Section 4 is devoted to estimating the model forecast error variance and its dynamics in time. Section 5 presents the conclusions and prospects for further developments.

## 2 The assimilation method and the numerical algorithm of its realization

### 2.1 The mathematical description of data assimilation scheme GKF

Let the mathematical model be governed by the equations

$$\frac{\partial X}{\partial t} = \Lambda(X, t) \tag{1}$$

with the initial condition $X(0) = X_0$.

Hereafter, $X$ means the model state vector defined on a phase-space, i.e. on the set of values which the model variables can take, $t$ is time, $\Lambda$ denotes the vector-function defined on the same phase-space and on a time-interval $[t_0, T]$. Without loss of generality $t_0$ is assumed to be 0. In the discrete realization, the model state vector has a dimension $r$, where $r$ is the number of grid points multiplied by the number of independent model variables. On the time interval $[0, T]$ the discretization $0 =$

$t_0 < t_1 < \ldots < t_N = T$ is introduced. For simplicity and also without loss of generality all these moments are assumed to be equidistant $\Delta t = t_{n+1} - t_n$. On each time interval $[t_n, t_{n+1}]$, $n = 0,1,\ldots,N-1$ model equations are numerically solved and instant $t_{n+1}$ data assimilation is performed by the formulae (Belyaev et al., 2018)

$$X_{a,n+1} = X_{b,n+1} + K_{n+1}(Y_{n+1} - HX_{b,n+1}), \tag{2}$$

$$5 \quad K_{n+1} = (\sigma_{n+1}^2)^{-1}(\Lambda_{n+1} - C_{n+1})(H\Lambda_{n+1})^T Q_{n+1}^{-1}, \tag{3}$$

$$\sigma_{n+1}^2 = (H\Lambda_{n+1})^T Q_{n+1}^{-1}(H\Lambda_{n+1}), \tag{4}$$

where $X_{a,n}, X_{b,n}$, $n = 0,1,\ldots,N$ are the model state vectors before and after assimilation, respectively, i.e., the analysis and background; $Y_n$ is the observation vector at the same instant of time that has a dimension $m$, where $m$ is the number of observation points multiplied by the number of independently observed variables. It is assumed that $X_{a,0} = X_{b,0} = X_0$ is the

10 known initial condition. Next, $K$ is the gain matrix (analogous to the Kalman gain matrix) with dimension $r \times m$, $H$ is the observational projection matrix with a dimension $m \times r$.

This projection interpolates the observed model variables from the model grid points onto the points of observations and simultaneously eliminates all redundant, i.e. not observed, variables. As usual., the superscript T denotes the transpose of a vector and/or a matrix. All those variables are considered on a time interval with a unit length.

15 Two variables in Eqs. (2)–(4) are assumed to be known, namely $C_n$ which is the trend, defined by the formula $C_{n+1} = \frac{E(\hat{Y}_{n+1} - X_{a,n})}{\Delta t}$. The physical sense of this vector is the defference between observations and previously constructed model state (analysis), where symbol $E$, ordinary, stands for the mathematical expectation or ensemble average, $\hat{Y}_{n+1}$ means the extrapolated observational vector, that is a vector which coincides with observations for observational phase-space but is extrapolated over on entire model phase-space. $\Lambda_{n+1} = \frac{X_{b,n+1} - X_{a,n}}{\Delta t}$, and $Q_{n+1} = E(\tilde{Y}_{n+1} - HX_{a,n})(\tilde{Y}_{n+1} - HX_{a,n})^T$, where

20 $\tilde{Y}_{n+1} = Y_{n+1} - HC_{n+1}$. Actually, $\tilde{Y}_{n+1}$ is the anomaly relatively the observational trend. According to definition $E(\tilde{Y}_{n+1} - HX_{a,n}) = 0$. The specific methods of constructing of the vectors $C$ and matrix $Q$ are presented in the next part.

This scheme with all necessary and sufficient conditions was introduced in Belyaev et al. (2018).It was also shown that this scheme generalizes the standard Kalman scheme method which follows from Eqs. (2)–(4), if $C_n = 0$ and $X_{a,n}$ coincides with the ensemble average.

25 **2.2 The numerical definition of known parameters**

As is seen from (2)–(4), this algorithm can be applied to an arbitrary numerical model with any physically reasonable initially conditions and it gives the output result (analysis) at any instant of time $t_n, n = 0,1,\ldots,N$. To provide its correct realization, it is necessary and sufficient to set up two aforementioned parameters, namely, the observational drift vector $C_n$ and the error covariant matrix $Q_n$. To set up the vector $C_n$, we apply the following algorithm: previously, using the Monte-

30 Carlo method the ensemble statistics from $M$ independent model runs is created. Theoretically, this ensemble should statistically represent the unknown "truth" value of field $X$ which satisfies the Eq. (1). Let this ensemble be denoted as $\hat{X}_n^j, n = 0,\ldots,N; j = 1,\ldots,M$ at each of moment of assimilation. Let us also assume that the analysis fields $X_{a,n}, n = 0,\ldots,k$, started from known initial condition $X_0$ to the instant of time $k$ where $k<N$ are already constructed. Then, the observational trend $C_{k+1}$ for each grid point is found out according to formula

$$35 \quad C_{n+1} = \left(M^{-1}\sum_{j=1}^{M}\hat{X}_{n+1}^j - X_{a,n}\right)/\Delta t. \tag{5}$$

The covariance $Q_n$ may be found in a similar way. If the ensemble statistics for moment of assimiation $t_i$ is known, then $Q_n$ is calculated by using the formula

$$Q_{n+1} = M^{-1} \sum_{j=1}^{M} (H\hat{X}_{n+1}^{j} - HX_{a,n})(H\hat{X}_{n+1}^{j} - HX_{a,n})^{\mathrm{T}}. \qquad (6)$$

It should be noted that in Eq. (6) we used only the observed components of the entire model state vector $X$. In the expanded form of Eq. (6) we consider the multiplication of the values at each pair of grid points but only between observed variables.

In Belyaev et al. (2018), it is proved that if the condition $E(\tilde{Y}_{n+1} - HX_{a,n}) = 0$ holds, than this construction of drift vector $C_n$ and covariance matrix $Q_n$ really estimates the observational trend and error covariance matrix of observational anomalies. Finally, using the Birkhoff –Khinchin theorem (Kolmogorov, 1938), which states that the ensemble average is approximated by the temporal average for sufficiently large numbers of series, we can rewrite formula (5) as follows:

$$C_{n+1} = \left( (n+1)^{-1} \sum_{i=1}^{n+1} X_{a,i} - X_{a,n} \right) / \Delta t, \qquad (7)$$

where $i$ is the number of assimilation time-steps until assimilation instant of time $n+1$. Our experiments have shown that for $n > 10$ formula (7) provides a good approximation of formula (5) ), however requiring much less amount of computations.


## 3 Computational experiments

### 3.1 The model and observational data base

Several numerical experiments were carried out, using the AVISO data and HYCOM model. The HYCOM model which was used only as a tool to perform the assimilation experiments is very well known (Bleck, 2002; Bleck et al., 1981) and there is no need to give its detailed description. The configuration of this model was as follows: its version 2.2.14 has a spatial resolution of approximately 0.25 degree in the Atlantic in both West-East (OX axis) directions with 420 grid points, and South-North (OY axis) with 720 grid points, respectively. This version considers 21 vertical layers of equal density within each of layer from the top to the bottom (OZ axis) The model domain covers the major part of the Atlantic from Antarctica and up to 55°N. It includes the Caribbean Sea and Mexican Bay, however excludes the Mediterranean Sea. On the lateral boundaries and on the sea bottom the Dirichlet conditions, i.e. the fixed climatological values are set up.

The model computes 4 barotropic variables (sea level, two velocity components and barotropic pressure on the sea surface) and 105 baroclinic components, namely, 5 variables on each given density layer: temperature, salinity, two velocity components, and layer thickness. After assimilation all of those variables change with respect to Eqs. (2)–(4). the total dimension of the model state vector denoted above as $r$ was 480x720x109.

For assimilation the AVISO data of sea level anomalies (SLA) have been used, where the temporal average for 10 years (2002-2011yr) has been used. This data has been downloaded and prepared for assimilation from website ([www.aviso.org](http://www.aviso.org)). The model data are interpolated onto observational points taken along satellite tracks. We assimilate only SLA. As the result of the quality control a part of observational array has been excluded from consideration. However, about 10000 daily recorded data remain to be used in the assimilation experiments. Therefore, the dimension of vector $Y$ denoted above as $m$ was about 10000 per day. Consequently, the size of the matrices used in Eqs. (2)–(4) has an order $\sim 10^9$ and to their computations the parallel algorithms have been used.

It is also reasonable to mention that the several studies related to the assimilation of AVISO data into HYCOM were presented earlier, in particular, Tanajura et al. (2015), where the assimilation of sea level anomalies has been performed used the EnOI method.

### 3.2 Parallelization of computations

Since this DA method is new, it is necessary to describe briefly the program realization of the presented scheme. The assimilation is realized as a separate program module DAM (Data Assimilation Module) in Fortran 95 with the use of MPI (Message Passing Interface) library. The calculated model variables are put in the DAM that operates independently on the

model blocks and utilizes another processing decomposition of the model domain. It should be noted that the size of 3D model arrays requires several GB. <span style="color:red">The high parallelization which means the explicit difference scheme is based due to the independency of GKF scheme for each observation points for Eqs. (2)–(4).</span> It should be noted that it is necessary to store a large data volume in the random-access memory with a memory limit of about 1.5--2 GB at each node. To optimize the assimilation processes, all MPI-nodes were assembled into several blocks with 8 elements in each of them. This significantly accelerates the computational process and requires about 40 minutes of computer elapsed time at 24 nodes instead of 4.5 hours at one node. The DAM module was realized on the HPC "Lomonosov 2" of the Lomonosov Moscow State University (Voevodin et al., 2012).

### 3.3 The algorithm of modelling

The assimilation experiments were performed as follows. Previously, the Spin Up run for the HYCOM was executed for 40 years from the rest and forced with the NCEP climatological reanalysis (Kalnay et al., 2002) of the wind stress and heat fluxes. Then the last 10 years of the modelled output were saved and daily recorded. This means that we have an archive with 10 completely defined fields of 109 computed model variables for each calendar day from January 1 to December 31 (the day February 29 was excluded from consideration).

The numerical experiments start from January 1 of 2010, forced by real wind stress and heat fluxes recorded on this day from GFS (global forecast system, NCEP) and lasts 1 month until January 31, 2010. Previously, the meteorological data were interpolated onto the model grid. The assimilations were executed daily, according to formulae (2)--(4). The information from daily recorded model outputs was used further to estimate the observed trend $C_n$ and observed anomalies $Q_n$ at the instant of assimilation n which coincides with the day number in the chronological order from 1 until 31.

The samples for the assimilation experiments according to formulae (5)--(7) were constructed as follows : the ensemble statistics for day $n$ were taken from the given climatological model calculation output with the same number and four other model outputs corresponding to two dates before and two dates after day $n$, equidistanly spaced with an interval of two days; so, in sum, there were 50 outputs. For example, to set up the ensemble statistics on January 15, we took the climatological model outputs on January 11, 13, 15, 17, and 19. Therefore, the sample length was 50.

We have carried out three types of experiments : A0 is the control experiment, where the model is integrated for one month without any assimilation; A1 is the assimilation experiment with the use of the EnOI method; and A02 is the assimilation experiments with the use of the GKF method.

### 3.4 Comparison with standard scheme EnOI

To compare the GKF method with the EnOI method, the standard EnOI scheme with the same data and the same model outputs used to create the ensemble was applied. Once it is done, the analysis is computed according to the Eq. (2), where

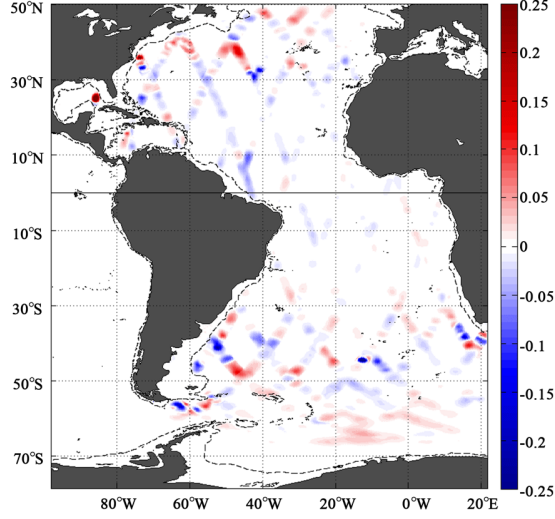$$K_{n+1} = \alpha B_{n+1} H^T (H B_{n+1} H^T + R)^{-1}, \tag{8}$$

$$B_{n+1} = M^{-1} \sum_{j=1}^{M} (\hat{X}_{n+1}^j - \bar{X}_{n+1})(\hat{X}_{n+1}^j - \bar{X}_{n+1})^T, \tag{9}$$

$\bar{X}_{n+1} = M^{-1} \sum_{j=1}^{M} \hat{X}_{n+1}^j$ means the ensemble average and for other notations we used the same ones as in Eqs. (2)–(4). The instrumental error covariant matrix $R$ and the empirical scalar $\alpha$ are defined "manually" from heuristic considerations. Formulae (8) and (9) are well known and, as was shown earlier in [16], they follow from Eqs. (2)–(4) if $C_{n+1} = 0$ and $\bar{X}_{n+1} = X_{a,n}$.

### 4 Results of the experiments and their analysis

<span style="color:blue">Sea level anomalies from AVISO as observational data have been used. There</span> data are recorded along the satellite

tracks (Fig. 1). On average, there are about 30 000 recorded values for each day, while the model computes much more sea level output values depending on the grid resolution. In order to make a correct comparison of the model outputs with observations, it is necessary to project the model values onto the observational points. The observed values may be both greater and smaller than the corresponding model values. In Fig 1, all the values for which the difference between the observations and projected model values is positive are indicated with the red color; otherwise, with the blue color.



**Figure 1.** Model domain and sea level anomalies (m) along the satellite tracks over the Atlantic. Red points show where data value exceeds the model one; blue points show the opposite.

To compare numerically the performance capabilities of different DA methods, we introduce the following quantities. Let

$$var_n = L^{-1} \sum_{i=1}^{L} ((SLA_m)_n^i - (SLA_o)_n^i)^2 \qquad (10)$$

be the variance of the model error at insant of time n, i.e., we square the difference between the model variable interpolated onto observational point i , in particularly, the sea level anomaly denoted $(SLA_m)_n^i$ and the corresponding observation value denoted $(SLA_o)_n^i$. This difference is taken over all observational points at each moment of assimilation $n$ from $n$=1,2,…,30 with the total number of observations $L$ . This number $L$ is diferent depending on the instant $n$ and it is not explicitly shown. Along with the variable $var_n$ we will consider two other variables $var_{f,n}$ and $var_{a,n}$ defined by formulae

$$var_{f,n} = L^{-1} \sum_{i=1}^{L} ((SLA_f)_n^i - (SLA_o)_n^i)^2, \qquad (11)$$
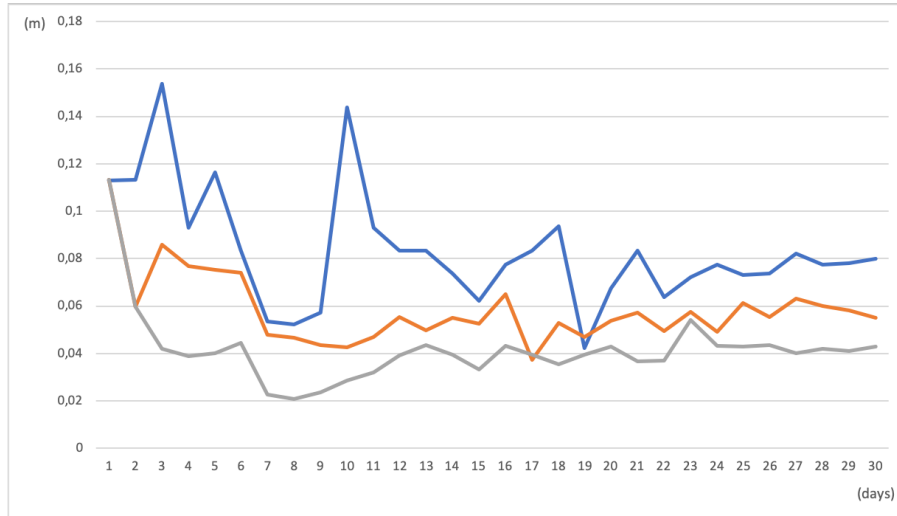
$$var_{a,n} = L^{-1} \sum_{i=1}^{L} ((SLA_a)_n^i - (SLA_o)_n^i)^2, \qquad (12)$$

where $(SLA_f)_n^i$ and $(SLA_a)_n^i$ are respectively the forecast and analysis model values at moment $n$ at observational point $i$. The forecast value $(SLA_f)_n^i$ means the model forecast at moment $n$ initialized from the analysis at the past time $n$–1, but the analysis value $(SLA_a)_n^i$ is counted at the same instant of time.

The time dependences of three root-mean-square variables $var_n$, $var_{f,n}$ and $var_{a,n}$ are presented in Fig. 2 and Fig. 3, respectively. Figure 2 contains the analysis error for three DA methods, including the control; Fig. 3 shows the 24h forecast errors for the same scheme. It is obvious that the control is unchanged in both Figs.

**Figure 2.** Analysis error variance of the SLA for 3 model runs. Blue line is the model control error, yellow line is the EnOI model error, gray line is the GKF model error.
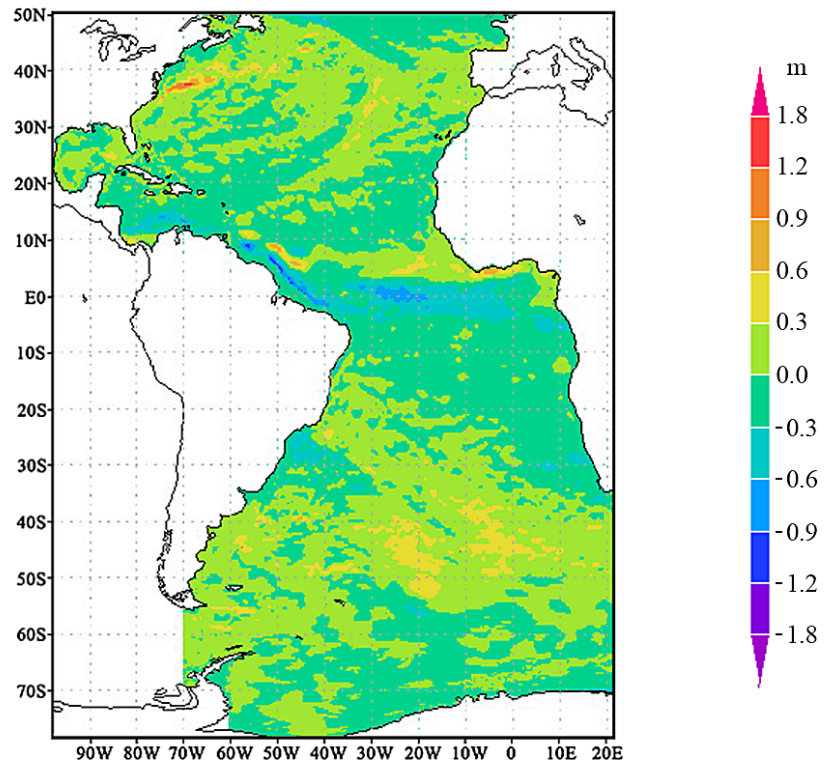


5 **Figure 3.** 24h forecast error variance of the SLA for 3 model runs. Blue line is the model control error, yellow line is the EnOI model error, gray line is the GKF model error.

As is seen from Figs. 2 and 3, the GKF scheme has a substantial advantage as compared with both the EnOI DA scheme and the control. Both the forecast and analysis errors in the GKF method are much smaller: one half--one third of the control error and one half of the EnOI error. We can also note that the EnOI scheme in general produces a smaller forecast error than

10 the control error; however, near day 20, this order is violated. For the same period, the forecast error of the GKF method is always smaller than the control error and nearly identical to the EnOI error near day 17. We can come to the similar conclusion about the analysis error. However, it should be pointed out that both DA methods operate properly; they both reduce the forecast error and improve the prognostic capacity of the model. We can also note that the model itself (control) makes the forecast error smaller with time since the model undergoes the sea surface temperature (SST) nudging

15 (relaxation), taken from the real NCEP/NCAR database. However, this procedure is beyond the scope of this paper.

Figures 2 and 3 show that the major discrepancy between both of DA methods and the control occurs near day 27, after which all the curves become practically steady. Therefore, our further analysis concerning model fields and independent data relates to this day, January 27, 2010.
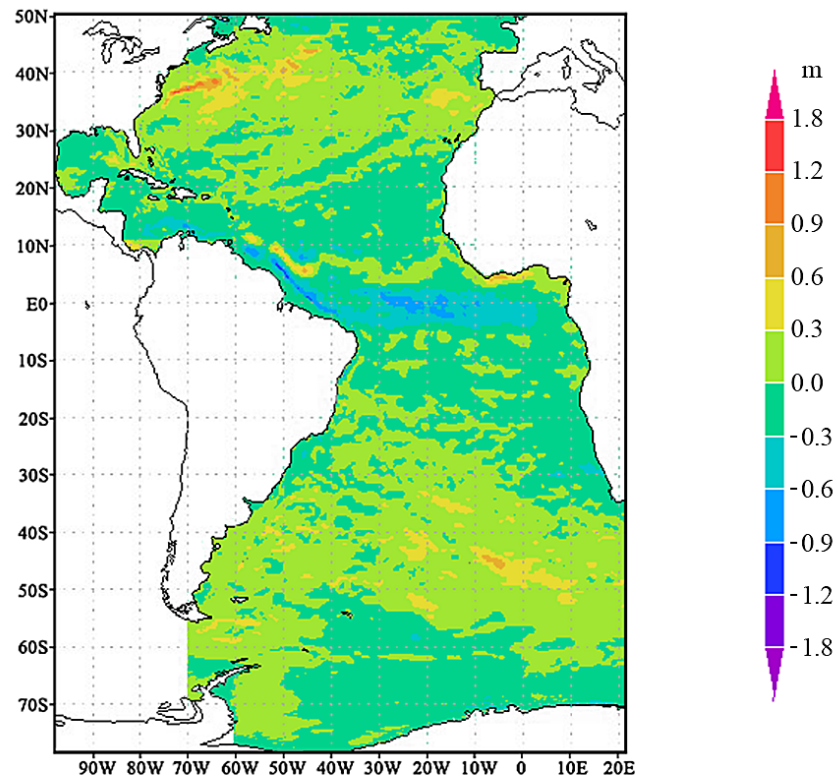
We start our analysis with the model output of SLA fields. Figures 4a--4d show the SLA model values, control, both

20 DA analysis, and their difference. Figures 5a--5d present the SST for these schemes as well.

All these figures show a similar SLA structure in the Atlantic with some specific details. Figure 4a (control) presents the general dynamic structure with pronounced positive eddies in the Gulf Stream zone and feebly marked negative eddies in the Brazil-Malvinas Confluence Zone in the Southern Atlantic. A negative SLA is also clearly seen at the equator; it propagates along the South America coastal zone in the Caribbean Sea and the Guinea Bay.
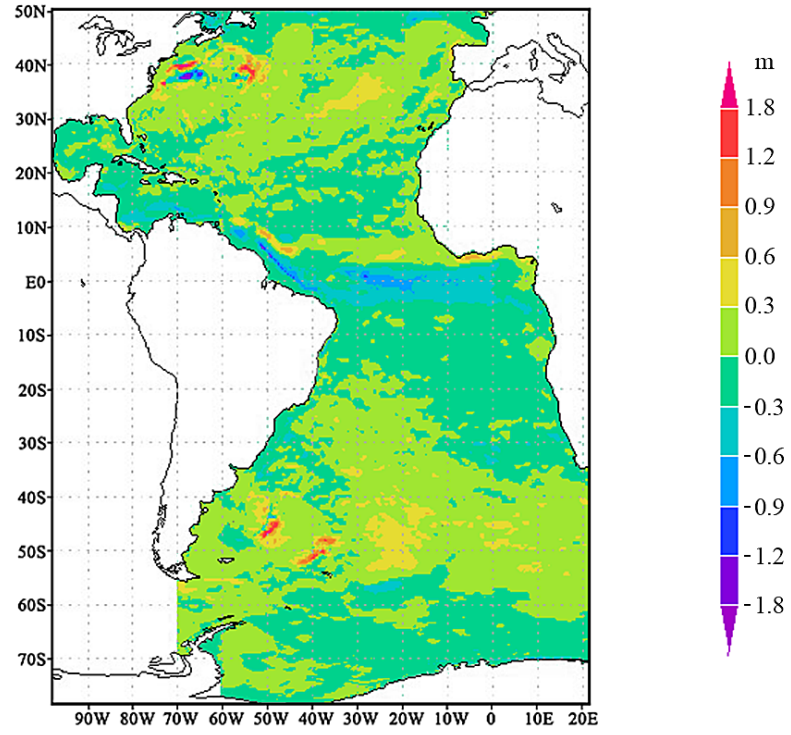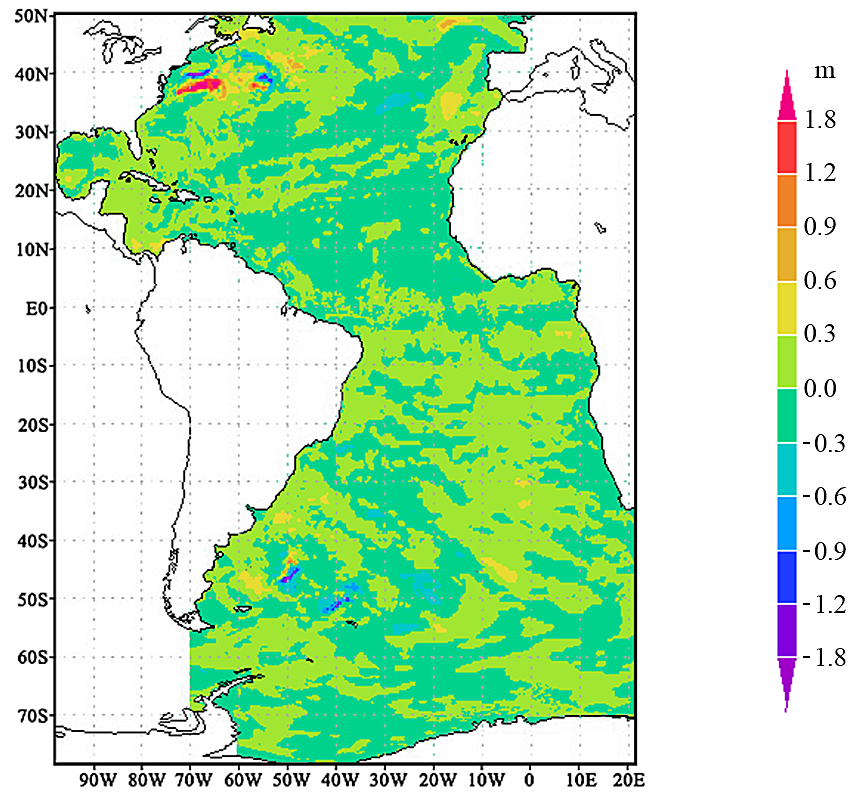
8

(a)



(b)

(c)



(d)



5    **Figure 4.** SLA after DA computation on 27.01.2010. (a) control, (b) EnOI, (c) GKF,(d) EnOI minus GKF.
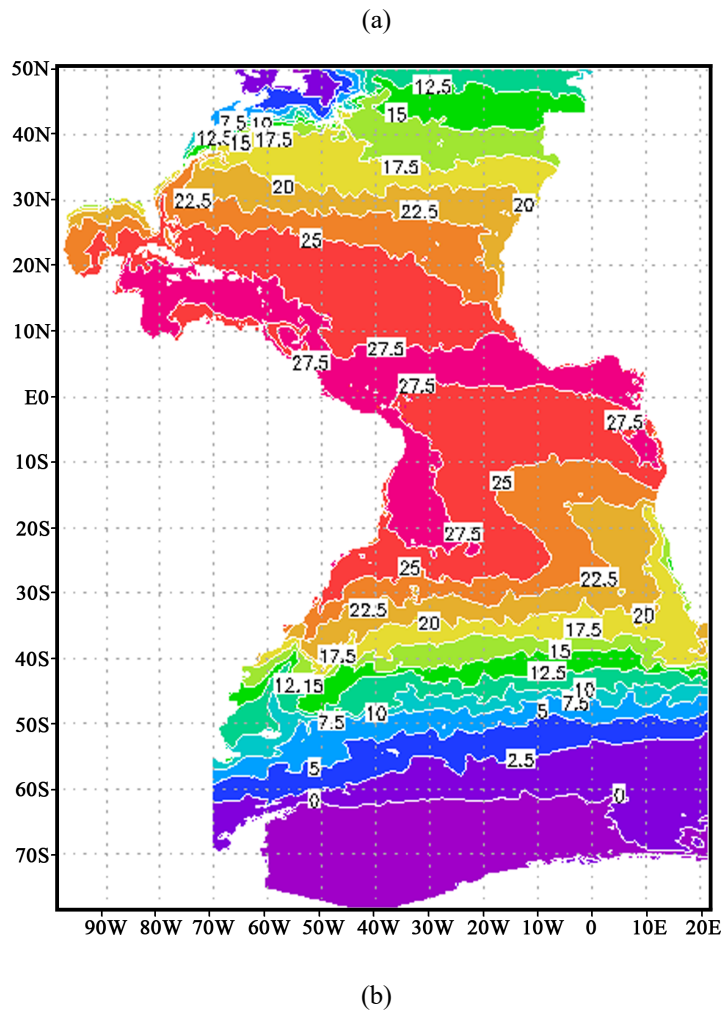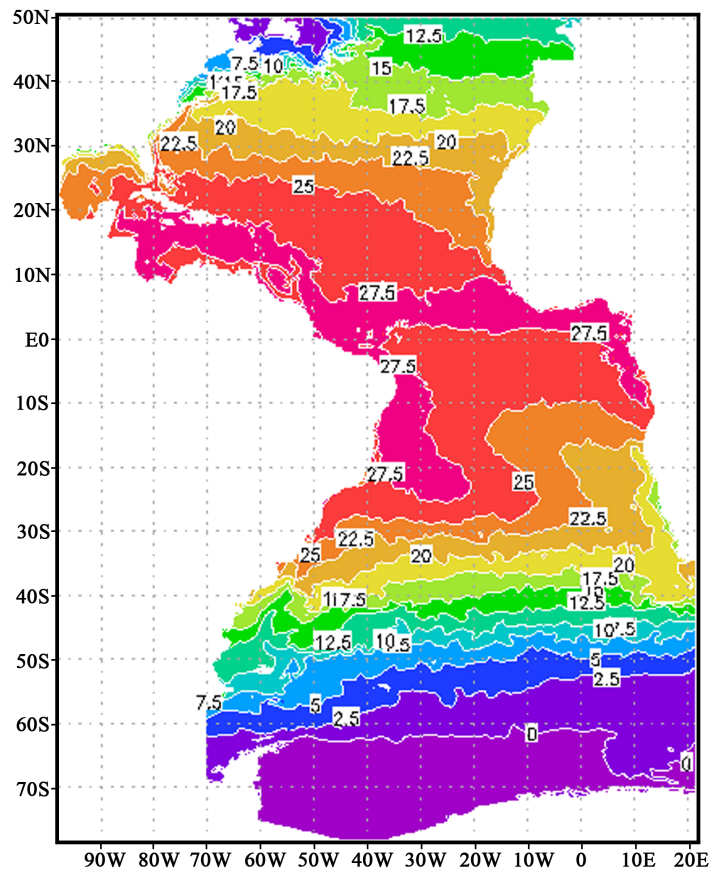
The EnOI (Kalman) filter DA method (Fig.4b), in general, reproduces the same structure; however, it essentially smooths the above-mentioned positive anomalies in the Gulf Stream zone and negative anomalies in the Center of the Atlantic. This can be explained by the fact that the EnOI filter method uses the ensemble statistics on 50 model outputs and they may considerably differ locally and instantly.

On the contrary, the GKF method substantially intensifies the positive anomalies in the Gulf Stream, makes them more compact and, at the same time, produces a positive anomaly in the Southern Atlantic near the Magellan passage. This is a temporal instant effect that corresponds to the local SLA trend in this zone. The difference between these two methods (EnOI minus GKF, Fig. 4d) confirms this conclusion; the positive anomaly appears neither in the control, nor in the EnOI scheme; however, it is present in the GKF method.
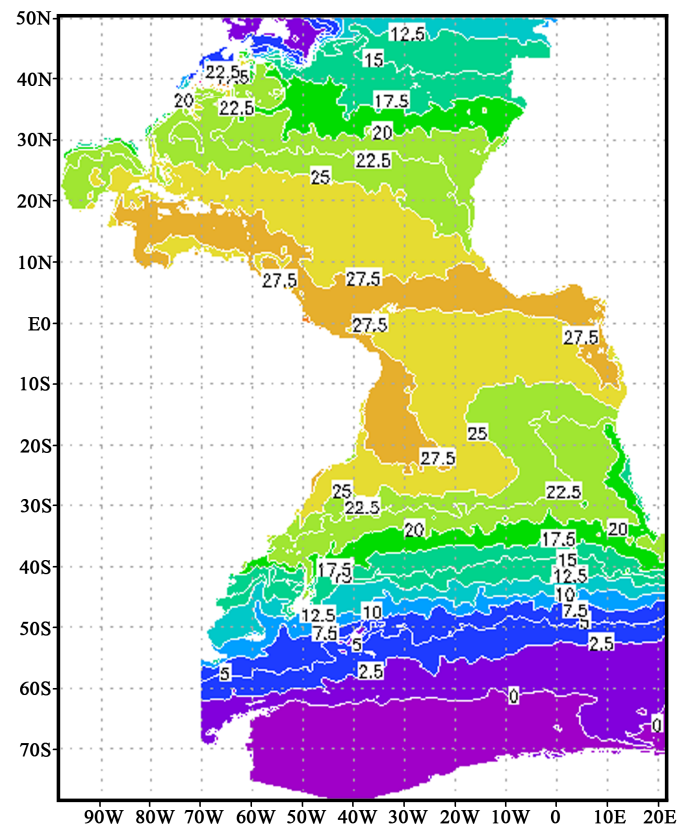
Figure 5 demonstrates the SST structure for the same domain and the same instant of time. As in Fig. 4, we can see in Fig. 5 the control SST (Fig. 5a), SST calculated with the use of the EnOI scheme (Fig. 5b), GKF scheme (Fig. 5c), and their difference (EnOI minus GKF). In fact, there is no visible difference between the control and EnOI calculations, some difference between them can be remarked in the northern part of the Gulf Stream zone. On the contrary, there is a considerable difference between the EnOI and GKF schemes, as is follows from Figs. 5c and 5d. This difference is clearly pronounced in the northern part of Atlantic, where a warm eddy appears and propagates along the current. This is a temporary effect and this meander is clearly pronounced locally. In the southern Atlantic, near the Brazil-Malvinas Confluence Zone, we can see a strong local dynamics. We can also assert that this is a temporary effect which is related to the instant time variability that is infinitesimal characteristic of the model vs data.
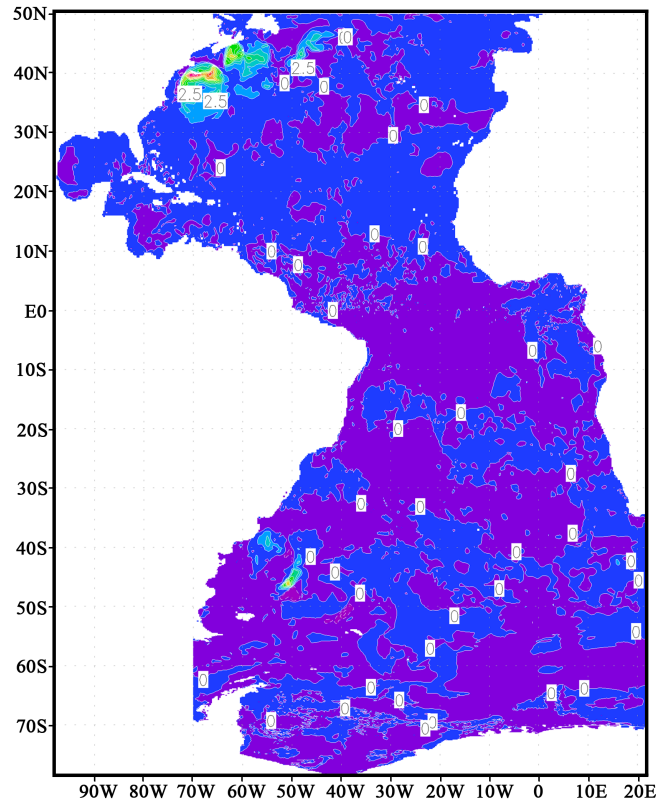
(a)



(b)

(c)



5

(d)

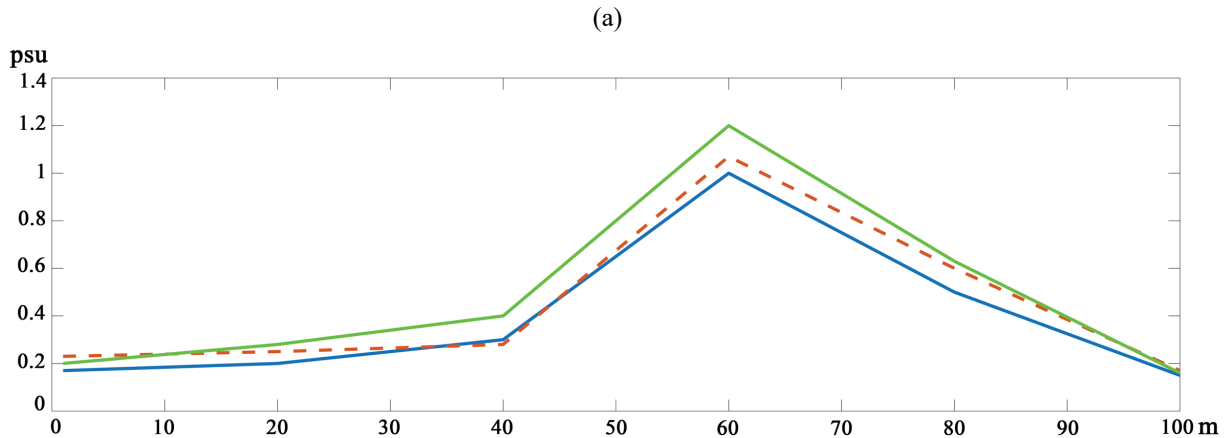**Figure 5**. SST for all assimilation methods. (a) control, (b) EnOI, (c) GKF, (d) EnOI minus GKF.
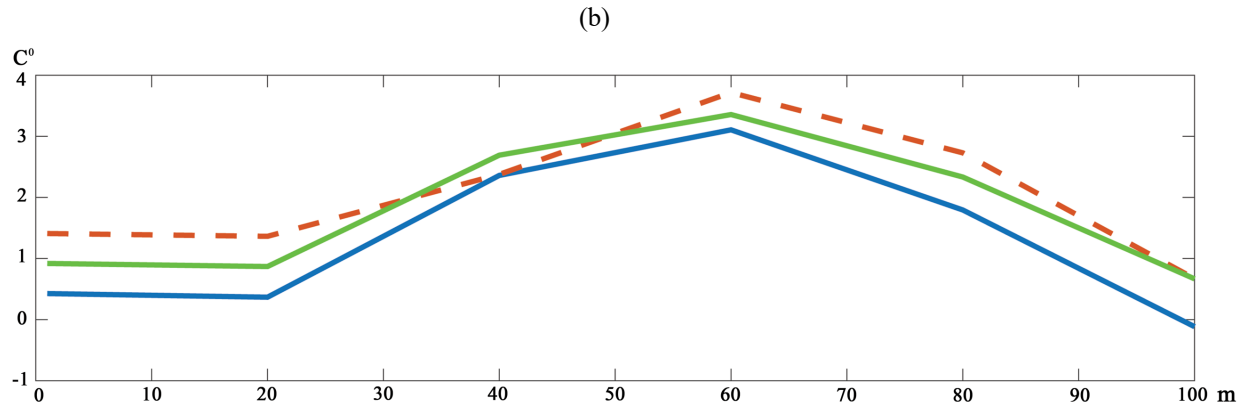
## 5 Comparison with independent data

### 5.1 Comparison with PIRATA moorings

For the comparison, we use the data array from the PIRATA moorins. As is well-known, there are 17 moored bouys in the Tropical and Center Atlantic; the temperature and salinity data from the sea surface to a depth of 500 m are recorded every 15 min and stored on the Internet. The daily data are availiable in the website http://pirata.ccst.inpe.br/en/data-2.

The comparison was realized as follows. On day 27 January, the average values from all model results with the two DA schemes and the control were linearly interoplated independently on each other onto the PIRATA bouys location and level. Then, the absolute difference between the model results and observations was calculated for all the bouys independently for each level. These computations were performed separately for the temperature and salinity.

The results of comparison are presented in Fig. 6. Figure 6a shows the deviation of the calculated salinity from the real data; Fig. 6b does the same for the temperature. The gray line refers to the control; the dashed line shows the deviation of the EnOI results from the real data; and the blue line corresponds to the deviation of the GKF results from the real data.

(a)

(b)



**Figure 6.** The difference with respect to the PIRATA mooring data: salinity (a), temperature (b).
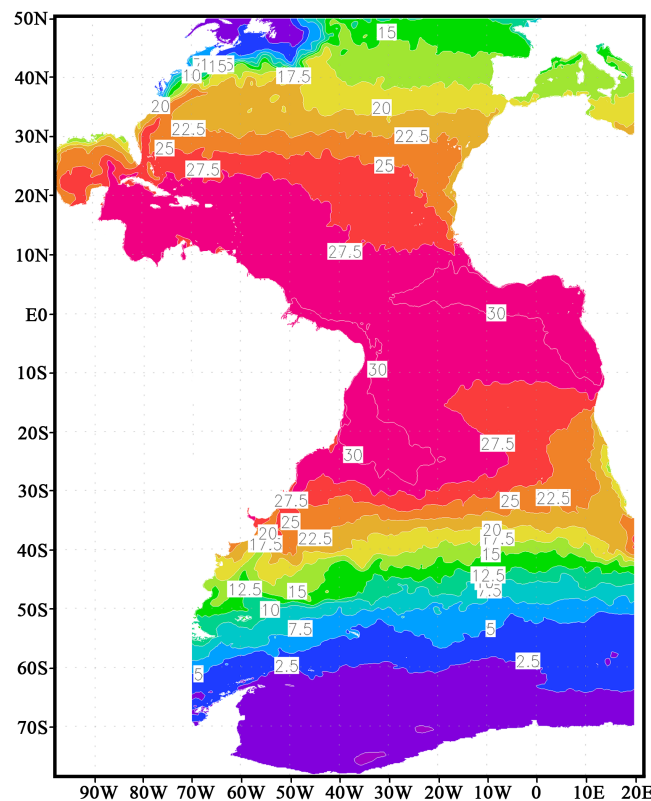
5　　As one can see from Fig. 6, the GKF method produces the best simulations in comparison with observations. According to Fig. 6a, the worst result for the salinity for all the methods, including the control, occurs at a depth of 40 m, however, for the temperature, it occurs at a depth of 60 m. Below 100 m, the temperature simulation is completely wrong and this can be explained by the fact that the HYCOM model does not describe adequately the equatorial area for the depths below 100 m. As for the  salinity,  the mixed layer is located exactly at a depth of 40 m and the model is unable to predict accurately its

10　position. However, all three schemes, including the control, are similar and reflect the reality.

**5.2 Comparison with SST observations**

　　Figure　　7　　contains　　the　　SST　　map　　directly　　downloaded　　from　　ftp://podaac-ftp.jpl.nasa.gov/allData/ghrsst/data/L4/GLOB/UKMO/OSTIA/2010/. Then, these data were interpolated onto the model grid. Since it is difficult to estimate all error sources, we do not show the numerical comparison of the model and observed SST;

15　we simply present the resulting SST field and qualitative comparison.



**Figure 7.** Observed SST field.

As we can note, the presented SST field contains all features characteristic for the Central and South Atlantic which are shown in Figs 5a--5c. However, Fig. 7 contains also visible synoptic eddies in the Gulf Stream zone, as well as in the Brazil-Malvinas Confluence Zone and these eddies are not seen in the control and EnOI computations (Figs. 5a and 5b). On the contrary, Fig. 5c contains pronounced eddies in both the Gulf Stream and Brazil-Malvinas zones, which agrees very well with the observed SST. At least, it is possible to assert that the GKF method is able to capture the synoptic variability in the ocean and does it better than the EnOI DA scheme and the control.

## 6 Conclusions

The novel feasible GKF DA method is presented and compared with the standard EnOI DA scheme which is often used in theoretical investigations and practical applications. The study proves both quantitatively and qualitatively that the presented GKF scheme has several advantages over the EnOI. In particular, it provides better the 24-h forecast and a posteriori analysis. In addition, the GKF method caputres the ocean synoptic-scale variability and its dynamics; the ocean temperature and salinity fields calculated after the application of the GKFmethod agree better with the independent PIRATA and OSTIA datasets. One of the advantages of the GKF method is its ability to calculate the confidence range of the forecast errors. We did not show it in this study but it was presented earlier in paper Belyaev et al. (2018). We have compared our new method GKF with EnOI because this two methods have the same both logical and methodological aproach. However, the GKF method may be compared with any other DA scheme with respect to the same criteria using the same model and dataset. In EnOI there are two parameters which allow one to control the results of DA, namely, scalar $\alpha$ and error observational matrix R. These two parameters (which can be, in fact, represented as the only parameter $R/\alpha$) are flexible and purely emphirical. They are set up manually and thear values may lie within wide intervals. However, GKF method operates with the same parameters since the matrices $HBH^T$ in EnOI and $Q$ in GKF, respectively, have the same physical meaning and play the same role in the assimilation scheme. The only - however, substantial – difference between GKF and EnOI consist in the fact that the $B$ matrix is calculated relatively to the model average state which is assumed to be equated to observational average, but in GKF it computes relatively to the linear trend (the difference between the current and previous model states) which is assumed to be equted to the average observational trend. It is weaker and more reasonable from the physical point of view, assumption becouse the reality can have a bias relative to the model, however this bias can be redused provided that the difference at two consequtive time moments is calculated.

The next studies in this direction can apply not only the SLA. but also SST, ARGO data as well as other available datasets. Similar studies can be performed for different regions and seasons to highlight the spatial and temporal variability of DA scheme.

*Data availability*. Data sets are available upon request by contacting the correspondence author.

*Author contributions*. The paper was written by KB and AK with the contribution by other authors. Numerical simulation was carried out by IS. Observational data and model integration have been provided and processed by CAST.

*Competing interests*. The authors declare that they have no conflict of interest.

# References

Agoshkov, V. I., Ipatova, V. M., Zalesnyi, V. B., Parmuzin, E. I., and Shutyaev V. P.: Problems of variational assimilation of observational data for ocean general circulation models and methods for their solution, Izvestiya, Atmospheric and Oceanic Physics, 46(6), 677–712, DOI: 10.1134/S0001433810060034, 2010.

5    Belyaev, K., Kuleshov, A., Tuchkova, N., and Tanajura, C. A. S.: An optimal data assimilation method and its application to the numerical simulation of the ocean dynamics, Mathematical and Computer Modelling of Dynamical Systems, 24(1), 12–25, https://doi.org/10.1080/13873954.2017.1338300, 2018.

Belyaev, K. P., Kuleshov, A. A., Smirnov, I. N., and Tanajura, C. A. S.: Comparison of data assimilation methods into hydrodynamic models of ocean circulation, Mathematical Models and Computer Simulations, 11(4), 564–574, DOI: 10    10.1134/S2070048219040045, 2019.

Belyaev, K. P., Tanajura, C. A. S., and Tuchkova, N. P.: Comparison of Argo drifter data assimilation methods for hydrodynamic models, Oceanology, 52(5), 523–615, DOI: 10.1134/S0001437012050025, 2012.

Bleck, R.: An oceanic general circulation model framed in hybrid isopycnic–Cartesian coordinates, Ocean Model., 4, 55–88, https://doi.org/10.1016/S1463-5003(01)00012-9, 2002.

15    Bleck, R., and Boudra, D. B.: Initial testing of a numerical ocean circulation model using a hybrid (quasi-isopycnic) vertical coordinate, J. Phys. Oceanogr., 11, 755–770, https://doi.org/10.1175/1520-0485(1981)011<0755:ITOANO>2.0.CO;2, 1981.

Cummings, J. A., and Smedstad, O. M.: Variational data assimilation for the global ocean, Data Assimilation for Atmospheric, Oceanic and Hydrologic Applications, Vol. II, Park, S. K., and Xu L. (Eds.), Springer-Verlag, Berlin, Heidelberg, 303–343, https://doi.org/10.1007/978-3-642-35088-7_13, 2013.

20    Evensen, G.: Data Assimilation – The Ensemble Kalman Filter, 2nd edn., Springer, Berlin, 2009.

Ghil, M., Malanotte-Rizzoli, P.: Data assimilation in meteorology and oceanography, Adv. Geophys., 33, 141–266, https://doi.org/10.1016/S0065-2687(08)60442-2, 1991.

Kalnay, E.: Atmospheric modeling, data assimilation and predictability, Cambridge University Press, New York, 2002.

Kaurkin, M. N., Ibraev, R. A., and Belyaev K. P.: Assimilation of the AVISO altimetry data into the Ocean dynamics model 25    with a high spatial resolution using Ensemble Optimal Interpolation (EnOI), Izvestiya, Atmospheric and Oceanic Physics, 54(1), 56–64, https://doi.org/10.1134/S0001433818010073, 2018.

Kolmogorov, A. N.: A simplified proof of the Birkhoff-Khinchin ergodic theorem, Uspekhi Mat. Nauk, 5, 52–56, 1938.

Lima M. O., Cirano, M.,. Mata, M. M., Goes, M., Goni, M., and Baringer, M.: An assessment of the Brazil current baroclinic structure and variability near 22° S in Distinct Ocean Forecasting and Analysis Systems, Ocean Dynamics, 66, 1–12, 30    https://doi.org/10.1007/s10236-016-0959-6, 2016.

Lorenc, A. C., Bowler, N. E., Clayton, A. M., Pring, S. R., and Fairbairn, D.: Comparison of Hybrid-4DEnVar and Hybrid-4DVar data assimilation methods for Global NWP. Mon. Wea. Rev., 143, 212–229, https://doi.org/10.1175/MWR-D-14-00195.1, 2015.

Marchuk, G. I., and Zalesny, V. B.: Modeling of the World Ocean circulation with the four-dimensional assimilation of 35    temperature and salinity fields, Izvestiya, Atmospheric and Oceanic Physics, 48(1), 15–29, https://link.springer.com/article/10.1134/S0001433812010070, 2012.

Penduff, T., Brasseur, P., Testut, C.-E., Barnier B., and Verron, J.: A four-year eddy-permitting assimilation of sea-surface temperature and altimetric data in the South Atlantic Ocean, Journal of Marine Research, 60(6), 805–833, https://doi.org/10.1357/002224002321505147, 2002.

40    Schiller, A., and Brassington, G. B. (Eds.): Operational Oceanography in the 21st Century, Springer, https://www.springer.com/gb/book/9789400703315, 2011.

Talagrand, O., and Courtier P.: Variational assimilation of meteorological observations with the adjoint vorticity equation I: Theory, J. Roy. Meteor. Soc., 113, 1311–1328, https://doi.org/10.1002/qj.49711347812, 1987.

Tanajura, C. A. S., and Belyaev K.: A sequential data assimilation method based on the properties of diffusion-type process, Appl. Math. Model., 33, 2165–2174, https://doi.org/10.1016/j.apm.2008.05.021, 2009.

Tanajura, C. A. S., and Lima, L.N.: Assimilation of sea surface height anomalies into HYCOM with an optimal interpolation scheme over the Atlantic Ocean METAREA V, Brazilian J. of Geophys., 31, 257–270,

5    http://dx.doi.org/10.22564/rbgf.v31i2.293, 2013.

Tanajura, C. A. S., Lima, L.N., and Belyaev, K. P.: Assmilation of sea level anomalies data into Hybrid Coordinate Ocean Model (HYCOM) over the Atlantic Ocean, Oceanology, 55(5), 667–678, https://doi.org/10.1134/S0001437015050161, 2015.

Van Leeuwen, P. J.: Representation errors and retrievals in linear and non-linear data assimilation, Quarterly Journal of the

10   Royal Meteorological Society, 141, 1612–1623, https://doi.org/10.1002/qj.2464, 2015.

Van Leeuwen, P. J.: Efficient nonlinear data-assimilation in geophysical fluid dynamics. Computers and Fluids, 46, 52–58, DOI: 10.1016/j.compfluid.2010.11.011, 2011.

Voevodin, Vl. V., Zhumatii, S.A., Sobolev, S. I., Antonov, A. S., Bryzgalov, P. A., Nikitenko, D. A., Stefanov, K. S., and Voevodin, Vad. V.: Praktika superkompiutera "Lomonosov", Otkrytye sistemy, Moscow, 7, 36–39, 2012. (in Russian)

15   Xie, J., and Zhu, J.: Ensemble optimal interpolation schemes for assimilating Argo profiles into a hybrid coordinate ocean model, Ocean Modelling, 33, 283–298, DOI: 10.1016/j.ocemod.2010.03.002, 2010.