



Forecasting the mixed-layer depth in the Northeast Atlantic: an ensemble approach, with uncertainties based on data from operational ocean forecasting systems

Y. Drillet¹, J. M. Lellouche¹, B. Levier¹, M. Drévilion¹, O. Le Galloudec¹, G. Reffray¹, C. Regnier¹, E. Greiner², and M. Clavier¹

¹Mercator Ocean, Toulouse, France

²CLS, Toulouse, France

Correspondence to: Y. Drillet (yann.drillet@mercator-ocean.fr)

Received: 28 April 2014 – Published in Ocean Sci. Discuss.: 11 June 2014

Revised: 30 October 2014 – Accepted: 17 November 2014 – Published: 12 December 2014

Abstract. Operational systems operated by Mercator Ocean provide daily ocean forecasts, and combining these forecasts we can produce ensemble forecast and uncertainty estimates. This study focuses on the mixed-layer depth in the Northeast Atlantic near the Porcupine Abyssal Plain for May 2013. This period is of interest for several reasons: (1) four Mercator Ocean operational systems provide daily forecasts at a horizontal resolution of 1/4, 1/12 and 1/36° with different physics; (2) glider deployment under the OSMOSIS project provides observation of the changes in mixed-layer depth; (3) the ocean stratifies in May, but mixing events induced by gale force wind are observed and forecast by the systems. Statistical scores and forecast error quantification for each system and for the combined products are presented. Skill scores indicate that forecasts are consistently better than persistence, and temporal correlations between forecast and observations are greater than 0.8 even for the 4-day forecast. The impact of atmospheric forecast error, and for the wind field in particular (miss or time delay of a wind burst forecast), is also quantified in terms of occurrence and intensity of mixing or stratification events.

Mersea, MyOcean and then MyOcean2 projects (<http://www.myocean.eu/>). Mercator Ocean is a French institution providing operational ocean forecasts for national requirements and also contributing to international efforts. A suite of global and regional ocean forecasting systems has been developed and provides daily forecasts. These forecasts are all available through MyOcean or Mercator Ocean services. All Mercator Ocean forecasting systems are evaluated in real time and in delayed mode with a well-defined protocol described in Lellouche et al. (2013). Results of this evaluation can be found in the quarterly “*QuO Va Dis?*” bulletin #13 for April–June 2013, available at <http://www.mercator-ocean.fr/eng/science/Qualification-validation2>. Part of the validation statistics are published on the MyOcean validation statistics webpage, available at <http://data.ncof.co.uk/calval/index.html>. Intercomparisons are also regularly performed in the context of the GODAE Oceanview international initiative (Ryan et al., 2014). These evaluations based on direct comparisons with observations (assimilated and independent) show that the four systems’ overall performance is state-of-the-art. The present study focuses on a “process oriented” validation criterion which is the mixed-layer depth at a specific location and time. An ensemble of mixed-layer depth estimates and the associated uncertainty are built around four operational forecasting systems. The Northeast Atlantic area was chosen because, since the launch of the V3 MyOcean service at the end of April 2013, four systems with different resolutions are providing ocean forecasts in this area on a daily basis. Moreover, glider observations for May 2013

1 Introduction

Operational oceanography has developed since the end of the 1990s in several countries with global-level partnerships under the GODAE Oceanview initiative (<https://www.godae-oceanview.org/>) and with European funding through

are available in the Coriolis database (available through My-Ocean service) sampling over the whole month with at least one profile per day in a small $1/2^\circ \times 1/2^\circ$ box centred on 16.25° W and 48.55° N. The physical variable chosen for this study is the mixed layer depth because this ocean variable is crucial in the forecast context for several reasons. (i) The mixed layer integrates a lot of physical ocean processes such as the horizontal and vertical advection and diffusion. Measuring the accuracy of the mixed-layer depth is an essential diagnostic for quantifying limitations in the model schemes or parameterizations (Giordani et al., 2005; Keerthi et al., 2013; Tozuka and Cronin, 2014). (ii) The mixed layer is the ocean layer which is directly in interaction with the atmosphere and the study of the mixed layer can reveal biases or unappropriated formulation of the atmospheric forcing of the ocean model (Béranger et al., 2010; Giordani, 2011). (iii) Most of the primary production in the ocean occurs in the mixed layer. Its evolution is partly driven by the horizontal advection in the mixed layer and by vertical processes at the base of the mixed layer which entrain nutrients from the deep layers into the surface layer (Lavigne et al., 2013). (iv) The mixed layer is directly in interaction with the atmosphere, and thus the vertical structure and heat content of this upper layer of the ocean is crucial in ocean–atmosphere coupling. This is well admitted in the scientific community and demonstrated at several scales for decadal (Meehl et al., 2014) or seasonal forecast (Balmaseda and Anderson, 2009), or short term forecast and especially for extreme events forecast as tropical cyclones (Goni and Trinanes, 2003). The ocean mixed-layer forecast is also important for defence applications such as estimating acoustic propagation or ambient noise (Shapiro et al., 2014). Our study aims at validating several ocean mixed-layer forecasts, and at better quantifying the influence of the various error sources, using ensemble techniques. This study focuses on a small region of the Northeast Atlantic in May, when the ocean stratifies and some mixing events occur which are directly linked to atmospheric forcing. The ability of a model to reproduce this critical stratification period validates its ability to reproduce the intraseasonal variability of the upper ocean. Moreover the spring stratification is a crucial phenomenon for the onset of phytoplankton blooms in this area (Mahadevan et al., 2012). Other studies quantifying the uncertainties in the ocean forecast for several oceanic fields (Lermusiaux et al., 2006) made use of super-ensemble techniques (Vandenbulcke et al., 2009; Lenartz et al., 2010; Pistoia, 2012; Scott et al., 2012) or have quantified the impact of medium-range atmospheric forecasting on the ocean (Drillet et al., 2009). An ensemble approach is also used in oceanography for estimating variability at a more climatic scale, for example in Zhu et al. (2012), Xue et al. (2012) and more recently in the Clivar Exchange special issue (<http://www.clivar.org/node/1507>). Some fairly complex techniques and diagnostics can be used, but in this study standard statistical techniques are used to compare several

estimates of the forecast mixed-layer depth. These forecasts come from several systems providing daily 5-day forecasts in our area of interest. Each ocean forecasting system includes a model with specific tunings, and oceanic initial conditions with specific data assimilation choices. Each system uses a slightly different set of atmospheric forcing extracted from real-time ECMWF atmospheric forecasts. Thus ensembles can be built from several model resolutions and tunings and several forecast lengths (time elapsed from ocean initialization), varying either the atmospheric forcing fields, or the type of ocean initial conditions. These two types of ensemble will help us quantify the impact of the atmospheric forcing errors, and the impact of the ocean initial condition errors on the accuracy of the mixed-layer forecast, depending on the forecast length. The paper is organized as follows: Sect. 2 describes the simulations and the observations used in the study. Section 3 draws on the statistics on the whole period, for quantifying forecast error and the main sources of forecast error. Section 4 describes the mixed-layer depth variability during May 2013, and how uncertainties in the observations and forecasts can be estimated. The last section presents the main conclusions of the study.

2 Forecast products and observations

The forecasts used in this study are provided by Mercator Ocean using four different operational systems. Two global ocean systems, one at $1/4^\circ$ horizontal resolution (Glo4, Lellouche et al., 2013) and the second at $1/12^\circ$ (Glo12), are used. Two regional systems are also used, one covering the North Atlantic and the Mediterranean at $1/12^\circ$ (Atl12, Lellouche et al., 2013) and the last at $1/36^\circ$ (Ibi36, Maraldi et al., 2013) covering the Northeast Atlantic. All these systems are based on the NEMO ocean code (Madec et al., 2008), using the same 50-level vertical grid and forced by ECMWF atmospheric analyses and forecasts. The initial state of each forecast is computed with data assimilation or with re-initialization techniques. The SAM2 software (Tranchant et al., 2008; Lellouche et al., 2013) is used to assimilate in situ and satellite observations. This reduced-order Kalman filter method is based on the singular evolutive extended Kalman filter (SEEK) formulation introduced by Pham et al. (1998). This method is used each week (on Wednesdays as shown in Fig. 1) to produce the initial state of the forecast for the Glo4, Glo12 and Atl12 systems. Two assimilation cycles are performed allowing the assimilation of observations up to 2 weeks old. The re-initialization method is used for the Ibi36 system (also on Wednesdays as shown in Fig. 1); a 2-week “spin up” is carried out to stabilize the high-resolution solution (at $1/36^\circ$) which is initialized using the $1/12^\circ$ analysis produced by the Atl12 system. This method and the effect of the length of spin-up time are detailed in Cailleau et al. (2012). The main characteristics of these systems are detailed in Table 1. Figure 1 shows more precisely how the

Table 1. Main characteristics of the ocean forecasting systems.

System	Glo4	Glo12	Atl12	Ibi36
Reference	PSY3QV3R3	PSY4QV2R2	PSY2QV4R4	IBI36QV2R1
Nemo	NEMO3.1			NEMO2.3 including specific development for regional/coastal application
Horizontal resolution	1/4° (~ 20 km)	1/12° (~ 6.5 km)	1/12° (~ 6.5 km)	1/36° (~ 2.2 km)
Vertical resolution	50 z vertical levels with partial step. 1 m at the surface. 22 levels in the upper 100 m.			
Atmospheric forcing	ECMWF operational analysis and forecast, spatial resolution ~ 12 km and 3 h temporal frequency. CORE Bulk formulation is used to compute atmospheric stress and fluxes.			
Atmospheric grid	Interpolated on 1/4° grid			Interpolated on 1/12° grid
Solar flux penetration	3-band parameterization for short-wave radiation (Lengaigne et al., 2007)			2-band parameterization for short-wave radiation (Morel et al., 2007)
Vertical mixing		TKE vertical mixing		GLS vertical mixing
Free surface		Filtered free surface		Explicit free surface with time splitting and tide
Initialization	SAM2 assimilation scheme (based on SEEK filter) assimilating SLA along track, L4 SST maps and in situ temperature and salinity profiles			Initialization with Atl12 analysis and 2-week spin-up.
Boundary conditions	none	none	from Glo4	from Atl12
Forecast length		7-day		5-day

systems are operated on a daily basis. Every day each system provides a hindcast estimate, H , of the ocean state. H is initialized using the “best” ocean state available, and forced with the “best” atmospheric forcing, i.e. the ECMWF analysis. The days following the simulation are F0, F1, F2, F3 and F4 – respectively the current day, the 1-day forecast and so on. All the ocean forecasts are forced by an atmospheric forecast. The latest available atmospheric forcing is used, and as the ocean forecasting systems are not launched at the same time they do not necessarily use the same 6-hourly update of the atmospheric forecast. Moreover, bulk formulae are used to compute atmospheric fluxes and this computation is another source of differences between the atmospheric forcing used by each system. Using this scenario, we can build a four-member ensemble for each forecast length differing mostly in their initial states, and the mean and median of this

ensemble can be considered as two other forecasts. We obtain for each date thirty 3-D ocean states which are not independent estimates of the ocean. A reference experiment, hereafter called Atl12 free, was also carried out using the Atl12 system without data assimilation. This reference was needed in order to assess the oceanic initial conditions errors (with respect to atmospheric forcing errors), and to discuss their impact on the mixed-layer depth forecast errors (Sect. 4.3.4). This simulation was initialized in March 2013 with the analysis provided by Atl12 system, and forced using the atmospheric forecast analysis to the end of May 2013.

This study focuses on May 2013, when all four ocean forecasting systems described above were providing forecasts and when the Coriolis in situ database contains repetitive in situ profiles obtained from glider observations in a 1/2° area centred on 16.25° W and 48.55° N (Fig. 2). These

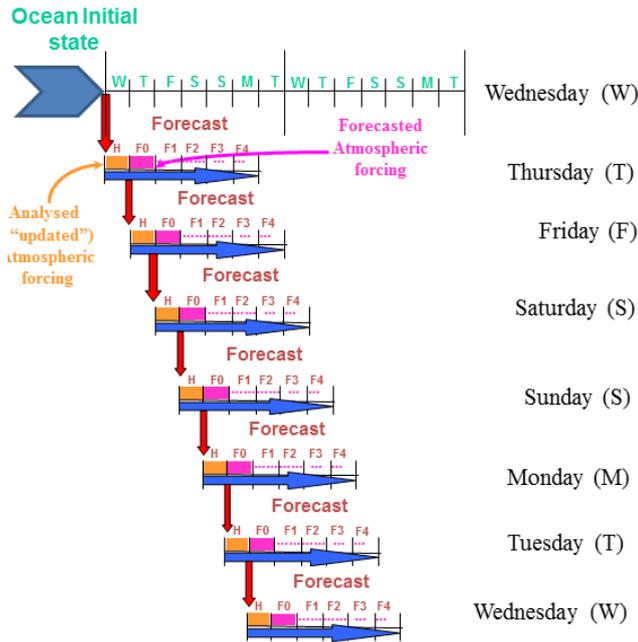


Figure 1. Operational scheme for producing daily forecasts with all the Mercator Ocean systems. The ocean initial state is produced once a week on Wednesdays. Then, starting from this state, a hind-cast (H) is produced each day using analysed atmospheric forcing. Then the forecast for the current day (F0) up to 4-day forecasts (F4) are performed daily, forced by the atmospheric forecasts.

gliders were deployed under the OSMOSIS (Ocean Surface Mixing, Ocean Sub-mesoscale Interaction Study, <http://www.bodc.ac.uk/projects/uk/osmosis/introduction/>) project, in which care is taken to apply a near-real-time quality control, disseminating the observations to the Coriolis data centre in real time. The observations used in this study do not represent the full data set but a subsampled one as carried out in all Mercator Ocean forecasting systems for all types of in situ observations. For each instrument (in this case gliders) only one profile per day is retained to avoid over-sampling (in time and space) of the observations since the global model cannot represent processes at such high resolution. This selection is made at random so that no particular time of the day is favoured. In our subsample database, 74 vertical profiles are available for May 2013 with at least one profile per day. This data set allows a good representation of the day-to-day temporal evolution of the temperature profile, and then of the mixed-layer depth during the entire month of May 2013. De Boyer Montégut et al. (2004) review temperature and density criteria which are used both for the model and the observations in computing the mixed-layer depth. In their paper, the criterion hereafter referred to as “0.2 °C temperature difference criterion” was used to compute global mixed-layer depth climatology based on in situ observations. In the present study the base of the mixed layer, for model and observations, as selected by this crite-

rior is the first level just above the point where the vertical temperature difference from the surface exceeds 0.2 °C. From a practical point of view, all the in situ profiles are interpolated on the vertical model grid to simplify the comparison between models and observations. Note that the precision of the mixed-layer depth estimate depends on the depth as function of the vertical grid; at the surface this is around 1 m and at 50 m depth around 10 m. Figure 3 shows the pertinence in our area of the “0.2 °C difference temperature criterion” which is also used in the routine validation of the operational production against in situ temperature profiles. A brief comparison with a density criterion is also performed in our area of interest (a 0.03 kg m^{-3} difference with the surface density is used to detect the base of the mixed layer, as suggested by De Boyer Montégut et al., 2004). When the mixed layer is really pronounced as for one profile on 11 May and for the two profiles of 28 May, the temperature criterion detects the base of the mixed layer whereas the density criterion gives a deeper mixed layer. When the profile is more mixed the base of the thermocline is also detected, but the density criterion indicates a shallower mixed layer (for instance in the homogeneous profile of 11 May) or similar values between the two criteria (for instance on 18 May when the three profiles are homogeneous). Each mixed-layer criterion gives a different estimate of mixed-layer depth for each individual profile. However, the mixed-layer depth time series computed with both criteria (not shown) display the same daily variability and amplitude throughout May 2013. This time period is of particular interest as it exhibits the spring re-stratification phase; gusts of wind occur also during this month and their effects on vertical mixing can be quantified. Some mesoscale oceanic structures are also present in this area, associated with strong fronts and currents which induce vertical mixing. In what follows, analyses and statistics computed using the model outputs and observations are based on (i) daily values which are actually daily means for the model but only the mean of all the data available during the day for the observations, and (ii) the spatial mean over the $1/2^\circ \times 1/2^\circ$ box defined previously. This box contains all in the situ profiles available during this month (Fig. 2) and is small enough when compared with the mesoscale structures in this area. This choice, both spatially and temporally, is justified by the fact that the model cannot simulate all the smaller scales available in the in situ observations. To illustrate more precisely the daily variability of the observations, Fig. 2 (right panel) shows a zoomed portion of selected dates. From 11 to 14 May there is a large variability in the observed Mixed Layer Depth (MLD) in a small $1/4^\circ \times 1/4^\circ$ box. All these observations should be in the same model grid cell in the $1/4^\circ$ model and in neighbouring grid cells in the $1/2^\circ$ models. However, observations show different profiles where the mixed-layer depth varies over several tens of metres as for example on 11 and 14 May. This cannot appear in the model because daily average outputs are stored. The best way to compare observations and model outputs at differing

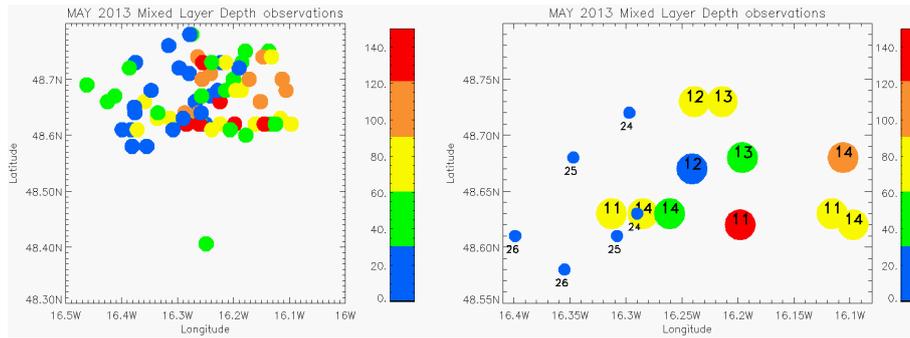


Figure 2. Left panel, position of the 74 profiles available in the area during May 2013. Right panel, selection of profiles from 11 to 14 May during the M1 mixing event and the first S1 re-stratification phase (large circles) and from 24 to 26 May during the S2 stratification phase (small circles). Colours show the mixed-layer depth computed for each profile with the 0.2 °C criterion. The number inside or below the circles gives the day of the measurement.

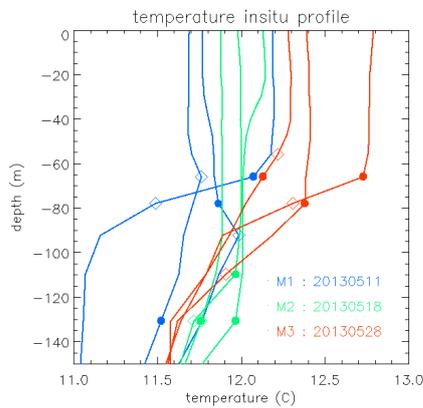


Figure 3. Available in situ profiles for three dates corresponding to three mixing events (M1, M2 and M3) during May 2013. Note that for these three dates three temperature profiles are available. The circles indicate the mixed-layer depth computed using the temperature profiles and the 0.2 °C criterion and the diamond using density profiles and the 0.03 kg m⁻³ criterion.

horizontal resolutions is to average spatially and temporally, and to consider daily profiles over the month where smaller scales in observation and high-resolution models are filtered out. Thus all estimates, both from models and from observations, will be compared on the same spatial and temporal grid. These daily profiles of a 1/2° × 1/2° box over a 1-month period is well suited for the physical mechanisms we study here. Note that by construction, nearly no spatial filtering is done when averaging Glo4 in 1/2° × 1/2° boxes, while the smallest scales are actually filtered out from Glo12, Atl12 and Ibi36.

3 Statistics

The statistics computed are mean bias (not shown), temporal correlation, error standard deviation for the Taylor dia-

gram (Fig. 4), skill score (Fig. 5) (Murphy, 1988) and root mean square error (RMSE; Table 2) for each system and for the ensemble mean and median. For each forecasting system and each forecast length, the skill score (SS) is computed as follows:

$$SS = 1 - \frac{\sum_{m=1}^M (y_m - o_m)^2}{\sum_{m=1}^M (p_m - o_m)^2} \quad (1)$$

On a given date m , y_m is the forecast value, p_m the persistence and o_m the observation. M is the total number of days in May 2013. We computed the skill score with two different type of persistence. The persistence p_m is either the persistence of the initial condition of the forecast, or the persistence of the last observation available. In other words, in the first case p_m is equal to the initial condition of the y_m forecast, and in the second case it is equal to the observation available on the initial day of the y_m forecast. Observation allows the use of the same “reference” state (in this case the observations) to compare different systems. Additionally the temporal evolution of the ensemble average and standard deviation is shown in Fig. 6 for mixed-layer depth, wind, heat flux and fresh water flux. The mean bias (not shown) is small for Atl12 and Ibi36 (less than 2 m up to 3 days of forecast length) and is greater in the two global systems with values greater than 5 m. The 4-day forecast has the same bias amplitude with all systems (around 5 m) but with a negative bias for Atl12 and Ibi36 and a positive bias for Glo4 and Glo12. Generally, there is a positive bias in the Glo4 and Glo12 mixed-layer depth, which means that the mixed-layer depth is too deep when it is underestimated with Ibi36. These results are consistent with the validation work done regularly for the Mercator Ocean real-time production (see “*QuO Va Dis?*” bulletins

available at www.mercator-ocean.fr/eng/science/Qualification-validation2). The Taylor diagram (Fig. 4)

Table 2. RMSE in metres for the mixed-layer depth computed with the systems, the mean value and the mean after removing one system. The F0, F2 and F4 forecast lengths are shown. For each forecast length the best forecast is bold underlined, and the other forecast with error not greater than 1 m compared with the best is shown in bold.

System	Ibi36	Atl12	Glo4	Glo12	Mean	M-Ibi36	M-Atl12	M-Glo4	M-Glo12
F0 RMSE	15.5	16.0	27.4	19.8	17.0	18.7	17.8	15.3	17.6
F2 RMSE	16.5	18.1	29.4	21.6	18.2	20.2	19.1	16.7	18.7
F4 RMSE	18.8	19.1	29.8	23.6	18.3	19.7	19.3	17.8	18.4

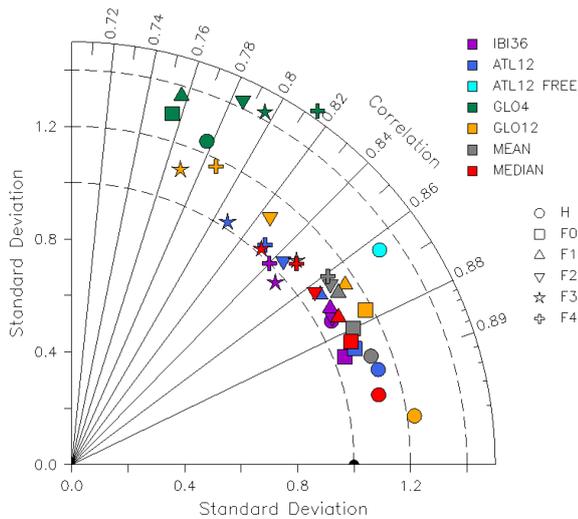


Figure 4. Taylor diagram comparing all available systems (in colour) and forecast lengths (symbol). The black dot with a standard deviation equal to 1 and a correlation of 1 indicates observations.

summarizes the following results: the temporal correlation between forecast and observation is greater than 0.85 for the first forecast day and decreases more or less depending on the system and/or the forecast length. Glo4 system is an exception; it has the lowest correlation for the first forecast lengths (from 0.78 for H to 0.76 for F0 and F1) and then increases to 0.81 for the 4-day forecast length. It can be diagnosed from Figs. 7 and 8 that the Glo4 F3 and F4 mixed-layer estimates have better results than the Glo4 H, F0 and F1 estimates for bad reasons. The poor correlation with observations of the Glo4 H, F0 and F1 mixed-layer estimates happens mainly because they miss the first stratification event S1 and the mixing event M2. Note that we will consider the ECMWF analysis, used to force the H estimates, as the observational reference for the winds. The S1 stratification takes place in the Glo4 F3 and F4 forecasts in response to underestimated winds from 6 to 7 May in F3, and from 6 to 10 May in F4. Then, the winds are stronger than observations in F3 and F4 just before the mixing event M2,

which induces mixing in Glo4 and improves the scores for F3 and F4 with respect to H, F0 and F1. The response to a realistic atmospheric forcing is not as good in

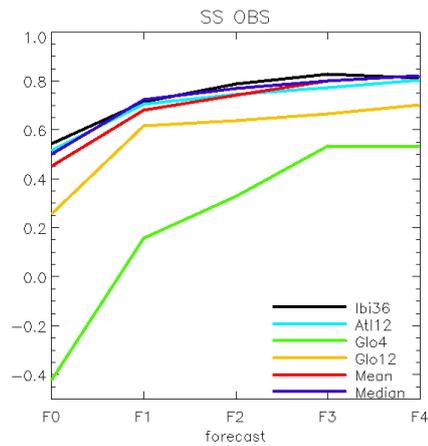


Figure 5. Skill score for the mixed-layer depth computed for all the systems and the ensemble mean and median during May 2013. The skill score is computed with the persistence of the observation

Glo4 as it is in Glo12, Atl12 or Ibi36, which will be discussed in more detail in Sects. 4.2 and 4.3. The ensemble mean gives the best result even if the Glo4 forecast is far worse than the other systems. One would expect the scores to decrease with the forecast length but the results are very similar (except for Glo4) for H up until the 1-day forecast; the dispersion of the systems (illustrated by the colour) is small in the Taylor diagram (Fig. 4) for all the metrics (correlation, standard deviation or rms). However, the forecast dispersion increases after the 2-day forecast and in particular there is a significant decrease in correlation to under 0.79 for Glo12, when it remains around 0.85 for Ibi36. The RMSE (Table 2) confirms previous results with a smaller error for Ibi36 and the ensemble mean (between 15 and 18 m rms) and a larger RMSE for Glo4 (between 27 and 30 m rms). The skill scores which measure improvement of the forecast in comparison to persistence of the initial condition (not shown) display very similar values as the one measuring the improvement in comparison to persistence of the last observation. The latter (Fig. 5) shows positive values (meaning that the forecast is better than persistence) for all forecasts except F0 in the Glo4 system. As expected, it increases with the forecast length meaning that the 4-day forecast is more efficient than the 1-day forecast in beating persistence. The largest skill score change lies between F0 and F1, meaning that after 1 day the forecast and

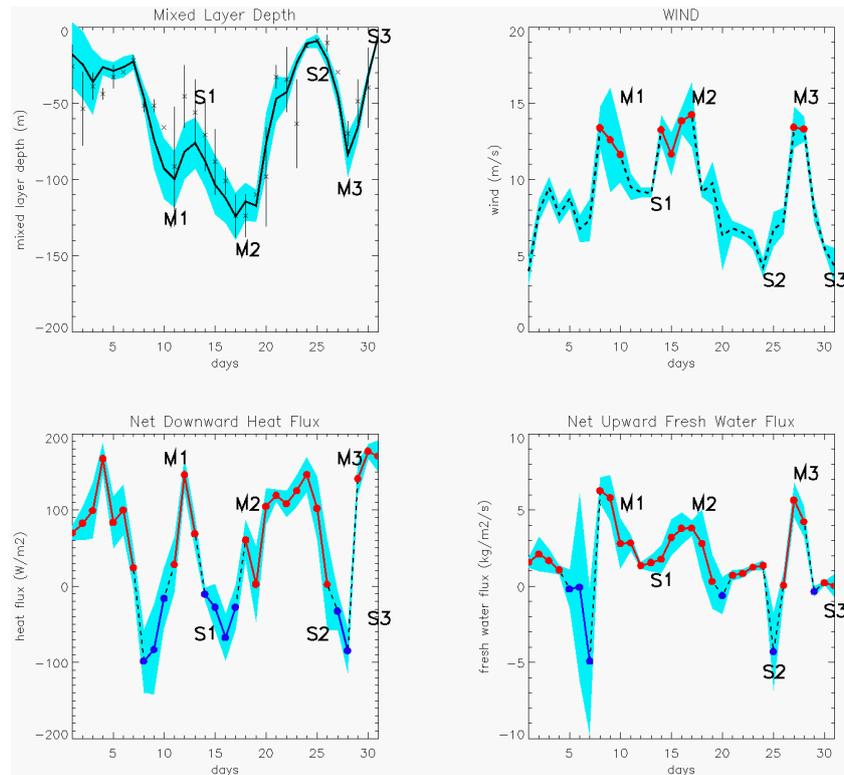


Figure 6. Top left: temporal evolution of the mixed layer simulated by the ensemble with the standard deviation in blue, and observations with associated uncertainties. Top right: wind speed time series. The analysis is in black dashed line with red line and dots when wind speed is larger than 9 m s^{-1} . The ± 1 standard deviation in blue is computed with all the forecast lengths; note that all systems are assumed to be using the same wind speed field, though an exception can occur if a forecast using one system is launched before atmospheric forcing is updated in the real-time production. Bottom left: total heat flux time series. The analysis is in black dashed line with blue curve and dots when the heat flux is negative (the ocean loses heat) and red when it is positive (the ocean get heat). The ± 1 standard deviation in light blue is computed with all forecast lengths and with all systems. Bottom right: fresh water flux time series. The analysis is in black dashed line with blue curve and dots when the fresh water flux is negative (precipitation term is dominant) and red when it is positive (evaporation term is dominant). The ± 1 standard deviation in light blue is computed with all forecast lengths and with all systems.

the last analysis or observation available are nearly equivalent, while after 2 days, the model has some significant predictive skill. Three “classes” of score can be seen as in the Taylor diagram (Fig. 4); the best is obtained with Atl12, Ibi36 and the mean and median products, a second with a significant decrease in the score obtained with Glo12, and a third with Glo4. Combining the forecasts in another way, simply by removing one system from the statistics, quantifies the gain (or degradation) obtained with each individual system. Table 2 shows the value of the RMSE for these combinations; the robust result is that the best forecast is obtained for the whole forecast with the mean computed after removing the Glo4 system, and with the Ibi36 system. Removing the Glo4 estimate, it may be noted that the mean of these forecasts is better than all the individual forecasts, showing that each estimate of the remaining ocean state gives pertinent information in terms of statistics for the forecast. In the following, the analysis of the mixing and stratification events will provide additional physical interpretation for these statistical results.

4 Mixed-layer depth forecast during May 2013

4.1 Description of the mixing and stratification events

In our area of study (Fig. 2), centred on 16.25° W – 48.55° N in the Northeast Atlantic, May 2013 is characterized by mixing and stratification events. Figure 6 illustrates this variability, with three mixing events (referred to as M1, M2 and M3) and three stratification events (referred to as S1, S2 and S3) well marked in both observations and simulations. Figure 6 shows the variability over the same period and in the same area for the main atmospheric forcing parameters, which are respectively wind speed, total downward heat flux and the upward fresh water budget. We note the good correspondence between the evolution of the mixed-layer depth and atmospheric forcing. The first mixing event (called M1 with a maximum value on 11 May) occurs just after a strong gust of wind ($\sim 13 \text{ m s}^{-1}$ on 8 May) and corresponds to an abrupt loss of heat which is negative from 8 to 10 May (minimal value $\sim 100 \text{ W m}^{-2}$ occurs on 8 May) and an evapo-

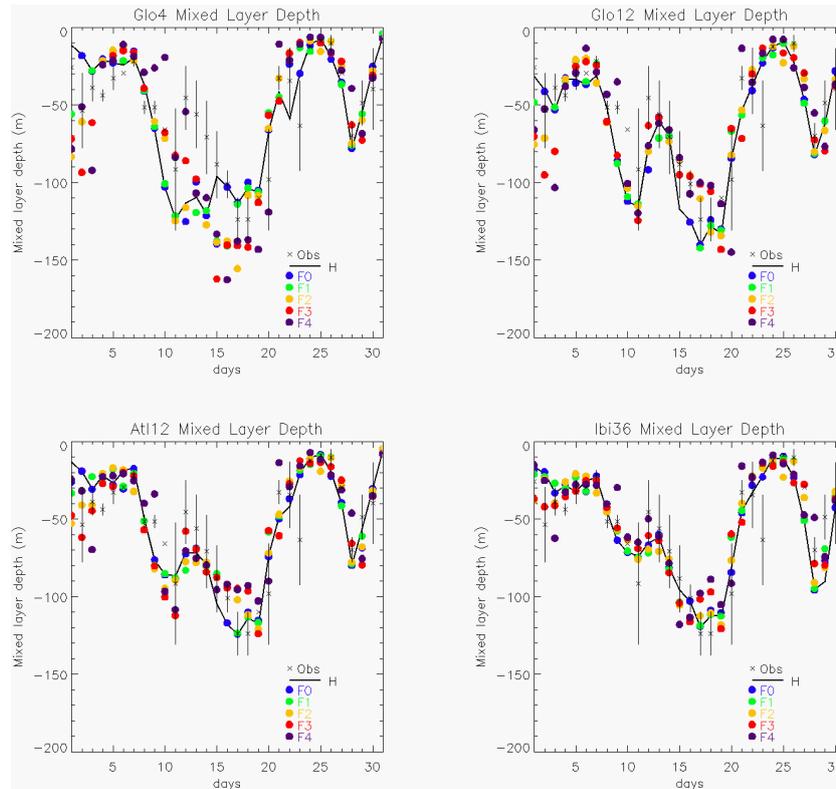


Figure 7. Mixed-layer depth evolution during May 2013. The black line is the hindcast and the coloured dots are the forecasts for several forecast lengths. The crosses are the means of the observations and the vertical black lines are error bars computed with the min and max values of the MLD estimated by the profiles during the day.

ration phase. The first stratification phase (called S1) is a short event occurring between 11 and 13 May and corresponds to a positive downward heat flux with a maximum value of $+150 \text{ W m}^{-2}$ on 12 May and a decrease of the wind to under 9 m s^{-1} . The second mixing event (called M2 with a maximum value on 17 May) is longer; it follows a short re-stratification phase before reaching the maximum mixed-layer depth and remains around 130 m depth for 3 days. This mixing phase is also preceded by strong winds and heat loss (negative downward heat flux from 14 to 17 May). A gradual stratification event (called S2) follows, occurring during a low wind and a warming period (from 18 to 26 May) which re-stratifies the entire water column. At the end of the month, a final strong gust of wind, causing heat loss, induces the M3 mixing event (28 May). The last rapid re-stratification of the entire water column (S3) occurs when the wind decreases. Several robust conclusions can be drawn from these alternating mixing and stratification events. First, all mixing events are associated with strong winds (exceeding 12 m s^{-1}) occurring a few days before the maximum of the mixed-layer depth is reached. For M1 and M2, the wind event occurs three days before the mixing maximum, while for M3 the response is faster (only 1 day). This could happen because the wind relative increase is larger (around 10 m s^{-1})

for M3 than for M1 and M2 (around $5\text{--}7 \text{ m s}^{-1}$) and because there have been 3 consecutive days of constantly increasing wind. These strong wind events are always associated with a large (less than -80 W m^{-2}) heat loss and evaporation. Re-stratification events occur when the wind speed decreases (to less than 9 m s^{-1}) and when the ocean absorbs heat with total fluxes greater than 100 W m^{-2} : 1 day for the S1 events and over a longer period (6 days) for the S2 event.

The standard deviation of all mixed-layer depths available for all systems and all forecast lengths computed in the same area centred on $16.25^\circ \text{ W}\text{--}48.55^\circ \text{ N}$ (Fig. 6) is also correlated with the uncertainties in the atmospheric fluxes, estimated as the standard deviation of all atmospheric flux estimates. There is a greater uncertainty for the mixed-layer depth during M1, S1, M2 and M3 with a standard deviation around 20 m, and also a smaller one around a few metres, during the S2 and S3 events. This is also true for the wind and heat fluxes where uncertainties are greatest during wind events, especially during the wind speed maxima. For the observations, the uncertainties are represented in Fig. 6 as vertical bars centred on the mean values of the observations in the box for every day in May 2013. Figure 2 shows the spatial distribution of these observations. The large uncertainty for the M1 and S1 events (from 11 to 14 May) is explained by the

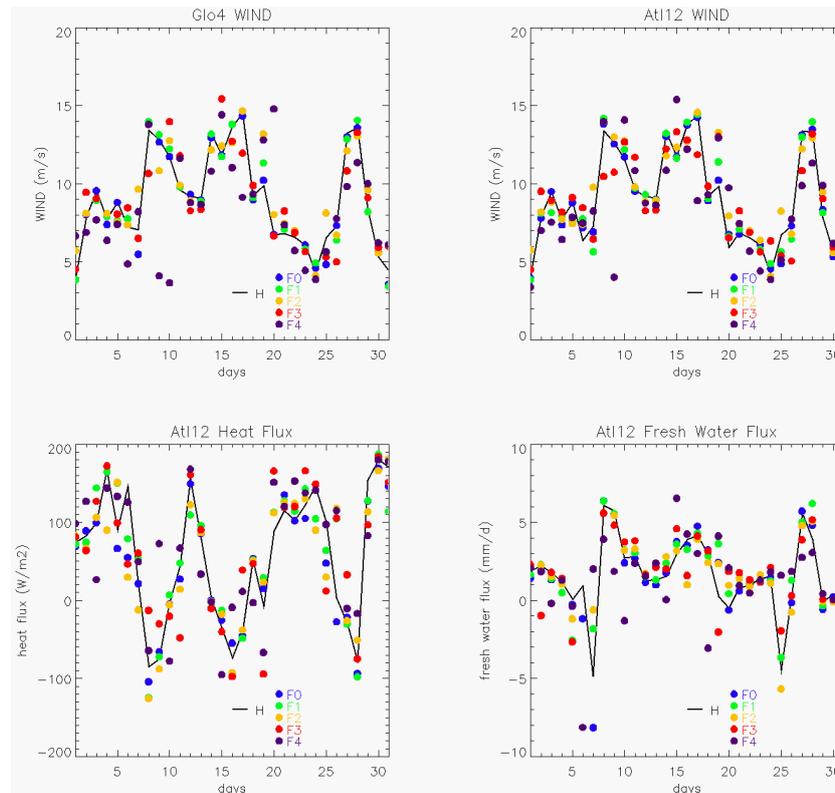


Figure 8. Temporal evolution of atmospheric forcing for hindcast (black line) and forecasts (coloured dots). Top panels: evolution of wind speed for Glo4 (left) and Atl12 (right) systems. Bottom panels: heat flux (left) and fresh water flux (right) for the Atl12 system.

fact that we have a large east–west and north–south gradient, with a centre area more stratified (with mixed layer around 40 m for 12, 13 and 14 May, the blue and green dots in Fig. 2) surrounded by more mixed profiles from 11 to 14 May (yellow, orange and red dots in Fig. 2). This gradient is smaller for the two other mixing events (M2 and M3). The uncertainty in the observations of the S2 stratification event is quite small and the right-hand panel of Fig. 2 shows a significantly shallow mixed layer of depth less than 30 m for 24, 25 and 26 May as indicated by the small blue circles. This uncertainty in the observation is not a robust diagnostic because the number of observations in our case is too small to give a precise estimate of this uncertainty, but nevertheless it gives useful information for evaluating the model. In this particular experiment, at this location, during this month and taking into account the estimate of the uncertainty for the model and the observations, the model is in agreement with the observations.

4.2 Evaluation of the hindcasts

Comparing the hindcasts (hereafter referred to as H) in Fig. 7 for the ocean fields and Fig. 8 for the atmospheric fields, all systems describe a stratified period at the beginning of the month with mixed-layer depth around 20 m, except for Glo12

where the mixed layer is deeper for the same period (around 40 m) which is closer to the observations. This may result from the large-scale conditions being different in Glo12 on the one side, and in Atl12 and Ibi36 on the other side. All systems have their own dynamical regime, but Ibi36 is initialized with Atl12 analyses, which explains that the circulation features of Atl12 and Ibi36 bear similarities, but do not look like the main circulation features of Glo12.

On 7 May all the systems simulate the beginning of the M1 mixing event, which reaches its maximum after 4 days but with significantly different amplitudes. The M1 event is too fast and too strong with Glo4 and Glo12 compared with the observations whereas the Ibi36 and Atl12 hindcasts are much closer to the observed values. Glo4 and Glo12 simulate mixed-layer depth greater than 100 m, while Atl12 simulates only 85 m of mixed-layer depth and Ibi36 even less so, with only 70 m depth. There is then a re-stratification event (S1), completely missed with Glo4, while it is observed and simulated with the other systems. The strongest re-stratification takes place with Glo12 while nothing happens with Glo4 where the mixed layer remains deeper than 100 m for 8 days. This stratification event is present in the observations, and the mixed-layer depths are very close to the observation in Glo12, Atl12 and Ibi36 even if the schedule of the re-stratification is different mainly due to differences

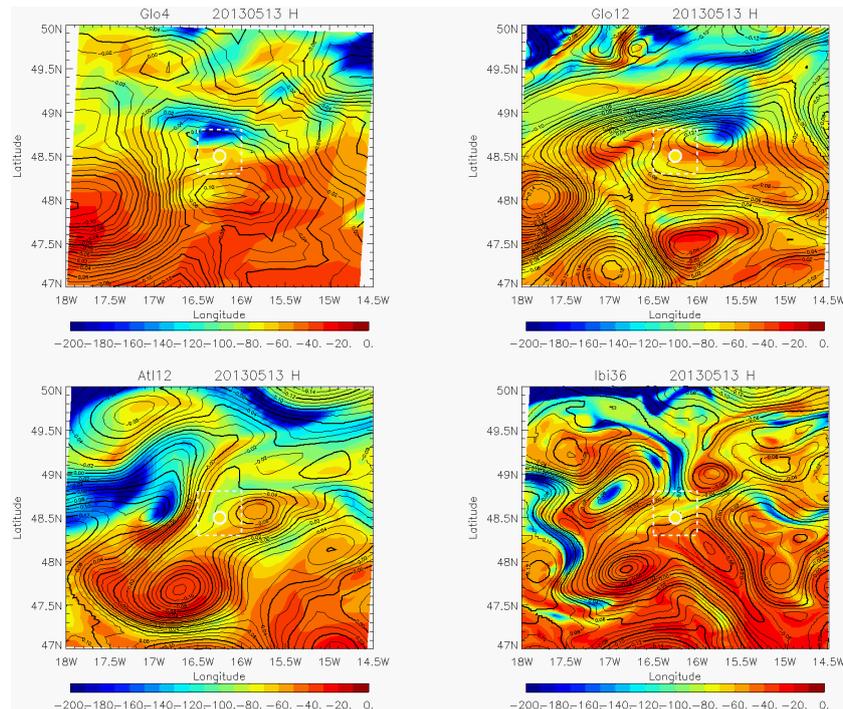


Figure 9. Mixed-layer depth (colour field) and sea surface height (contours) simulated by the four systems for 13 May in the area surrounding the area of interest ($1/2^\circ \times 1/2^\circ$ white box centred on 16.25° W– 48.55° N, Fig. 2). The colour dot indicates the mixed-layer depth observation for the same day.

during the first mixing event as noted above. Figures 9 and 10 show the spatial patterns of the mixed-layer depth for all systems for 13 May and 16 May respectively. In our area of interest (white squares in these figures) there is a strong gradient in the mixed-layer depth with a mixed column in the northern part of the area, and a more stratified ocean in the south. In this case the mean profile in this box is not fully representative of the situation and the observation fails to capture this kind of pattern. Nevertheless, as Figs. 9 and 10 show, the hindcast mixed layer in this box fit well with the mean observed mixed-layer depth. Statistics computed over a smaller box (taking into account only the northern part of the box from 48.55 to 48.8° N) are slightly different for the Glo4 system with a deeper M1 mixing event and a more stratified S1 event (not shown). But in this case the number of points in the box is too small for this low-resolution system, and the statistical results in terms of bias or rms values are not as good. As explained in Sect. 2, the average applied over the $1/2^\circ \times 1/2^\circ$ box is a small-scale filtering which is efficient for the $1/12^\circ$ or the $1/36^\circ$ of degree system and consistent with the available observations, but filters no signal for the $1/4^\circ$ system. Taking into account a larger box for this system could be a solution, but in this case the inconsistency with the available observations which are really concentrated in this small area will induce other biases. The M2 event with a maximum of mixed-layer depth on 17 and 18 May is well simulated with the Glo12, Atl12 and Ibi36 systems. The last

period of the month is more similar in all systems, with a re-stratification of the entire water column (S2) from 20 to 25 May, and a new mixing event (M3) followed by a re-stratification (S3). The temporal evolution of the mixed-layer depth agrees well among all the systems with minima and maxima occurring on the same day except for the S1 stratification event in Glo4 between 11 and 13 May. Observations available at this position allow a precise validation of the evolution of the mixed layer during the month. As shown by the statistics, the Ibi36 system is the closest to observations with very good timing of mixing and re-stratification events and a good estimate of the mixed-layer depth. This is not surprising as this system benefits from the highest-resolution horizontal mesh, the highest-resolution atmospheric forcing, and also from the GLS scheme for the vertical mixing (Umlauf and Burchard, 2003; Refray et al., 2014). For instance, on 13 May (Fig. 9) at 17° W a dipole structure can be detected in both Atl12 and Ibi36, with a cyclonic eddy near 48° N and an anticyclonic eddy near 49° N. This dipole structure is not present in either Glo4 or Glo12, which do not represent the anticyclone at 49° N.

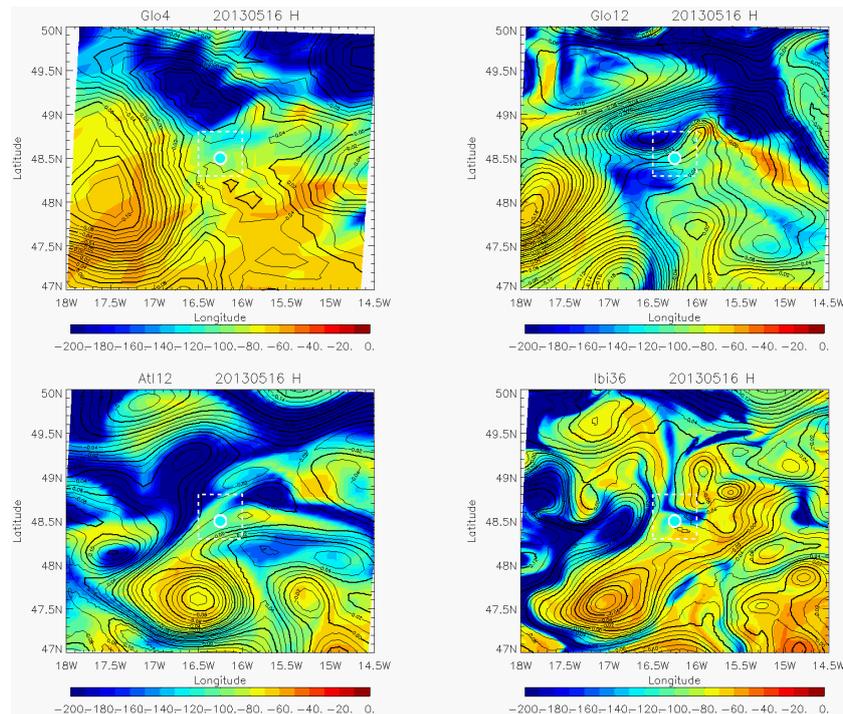


Figure 10. As Fig. 9, for 16 May.

4.3 Discussion of the forecasts

4.3.1 Forecast of the first mixing event (M1)

The greatest forecast error is obtained with the Glo4 system during the M1 event. During this first period (between 9 and 12 May) the 1- and 2-day forecasts are consistent with the hindcast (green and blue dots with respect to the black line in Fig. 7) and so deeper than the observations, but the 3- and 4-day forecasts (red and purple dots in Fig. 7) are closer to the observations with a thinner mixed layer. One would have expected F0 to be more accurate than F4. As already mentioned in Sect. 3, the Glo4 is better at long forecast lengths than at short forecast lengths for wrong reasons. The 4-day wind forecast is less than the analysis wind (4 m s^{-1} rather than 13 m s^{-1} ; purple dots for 9 May in Fig. 8, top panels). The Glo4 system seems to produce too much mixing in response to realistic wind (H, F0 to F3), and thus produces a more realistic mixed-layer depth when the wind is underestimated (F4). The other systems, Glo12, Atl12 and especially Ibi36 have a more realistic answer to wind forcing and for these systems the error increases with forecast length. At the beginning of the M1 event, the 4-day forecast misses the mixing. Looking more closely at the forecast for the 9 May (Fig. 7) none of the models 4-day forecast (purple dots) simulates the mixing when the smaller forecast lengths (blue, green, yellow and red dots) capture this event. The same kind of underestimation of the wind fields is observed for the Glo4 forecast at other dates and other forecast lengths (like 6 and

10 May for the 4-day length illustrated by the purple dots on Fig. 8, top left panel). The atmospheric forcings seen by Glo4 display slightly more dispersion between the different forecast lengths than the other systems. This can add to the uncertainty of the mixed-layer depth in Glo4 with respect to the other systems. These differences in the wind field used for the forecast are explained by the fact that in the operational suite all the systems are not launched at the same time. It is then possible to use the different base times of the atmospheric forecast for the ocean forecasts provided by the different systems used in this study (Glo12, Glo4, Atl12 and Ibi36). As the Glo4 system is the first to be launched in the operational suite, if there is a delay in the atmospheric forcing construction procedure, this system will use the latest atmospheric forecast (using for example the previous analysis cycle). The other systems are able to forecast this mixing of the water column up to 4 days. Glo12 and Atl12 provide an excess of mixing especially for the 3- and 4-day forecasts. Ibi36 is in better agreement with observations except for the 11 May where the observed mixed layer is deeper (a depth of 90 m but with high uncertainty) and the forecast, just as with the hindcast, gives too shallow a mixed layer (a depth of between 65 and 75 m).

4.3.2 Forecast of the first re-stratification (S1) and second mixing (M2) events

As already discussed for the hindcast in Sect. 4.2, the S1 re-stratification event is not forecast with Glo4. Although the

3- and 4-day forecast seem to give good results, it is for the wrong reason; the initial state of these forecasts is too stratified and the strong wind event is not present in the atmospheric forecast. The other systems are able to forecast this re-stratification phase after the 12 May for each forecast length. During the second mixing event (from 12 to 17 May in the observation) the Glo4 forecast (especially from day 2 to day 4) provides a deep mixed layer, deeper than the hindcast and also deeper than the observations. The analysis of the area (Figs. 9 and 10) shows that all systems provide mixing of the water column from 13 to 16 May. This is true for the hindcasts (Figs. 9 and 10) and forecasts (not shown) but at a larger scale than the smaller $1/2^\circ \times 1/2^\circ$ box which contains the observations, and which is illustrated by the white box in the figures. At this small scale, mesoscale oceanic structures affect the mixed layer and a new source of uncertainty is added to the atmospheric forcing uncertainties. As observed in Figs. 9 and 10, similar large-scale mixed-layer depth patterns appear in all systems, with a north–south gradient with shallow mixed layer in the south (less than 50 m depth) and a deeper mixed layer in the northern part. Note that the figures show hindcast states and consequently the atmospheric uncertainty is reduced. At smaller scales, the effects of mesoscale, fronts, eddies and associated dynamics are represented by the contours of sea surface height in Figs. 9 and 10. In this case it is noticeable that the horizontal resolution of the system is a key factor in the effect on the mixed-layer depth. In Glo4, at $1/4^\circ$ resolution, there is less consistency between the mixed layer and the sea surface height fields; at $1/12^\circ$ (in Glo12 and Atl12) and even more so at $1/36^\circ$ (Ibi36) there are thin structures along fronts, surrounding eddies where the mixed layer is deeper. This influences the statistics when looking at small spatial and temporal scales, as in our case where the spatial scale is less than 50 km and the temporal scale is approximately 1 day. As mentioned in previous sections, this S1 to M2 period contains uncertainties for the mixed-layer depth and also for the atmospheric forcing. It is linked to the following phenomena, which all contain uncertainties:

1. error in the atmospheric forecast (see Fig. 8);
2. rapid stratification/mixing change occurring over two days; in this case a short delay in the forecast gives a large error;
3. M2 event occurs when the mixed layer is still thick; in the case of a shallow mixed layer, the uncertainty is naturally reduced;
4. there are well marked mesoscale structures which affect the mixed-layer depth, generating vertical mixing associated with vertical velocities along the front and around eddies.

4.3.3 Forecast of the second and third stratification (S2, S3) and third mixing (M3) events

The S2, M3, S3 time sequence is well forecast in all the systems, with good temporal consistency with observations (Fig. 7). Maximum stratification occurs on 25 May (S2). Then, the water column is mixed until 28 May (M3) and quickly re-stratified until the end of the month (S3). All the forecast lengths are close to the hindcast run except the 4-day forecast for 21 and 28 May. For these dates, all systems give consistent solutions with too rapid a re-stratification for 21 May and a lack of mixing for 28 May. This is explained by the error in the wind forecast (Fig. 8) taking into account a 1- or 2-day lag, which is the typical time taken to mix the water column. For 19 and 20 May the forecast wind speed is too strong with wind speeds exceeding 10 m s^{-1} , while analyses give values less than 10 m s^{-1} decreasing to 7 m s^{-1} for 20 May. The opposite occurs for 27 May with a wind forecast of approximately 10 m s^{-1} rather than the 14 m s^{-1} predicted by the analysis

4.3.4 Atmospheric forcing versus initial state in the uncertainties

The question of the significance or effect of atmospheric forcing vs. initial state on the mixed-layer forecast has to be addressed. One diagnostic computed to quantify these two aspects separately is based on the temporal correlation between several time series. The first step is to compute the temporal correlation between the same forecast lengths with all the available systems. Correlations are thus computed for six ensembles of estimates (H, F0, F1, F2, F3, F4), each ensemble being made of four different time series coming from the four systems. If the initial state had a strong impact on error growth, one would expect the mean correlation of the ensembles to decrease significantly with the forecast length. In this case the mean correlation decreases from 0.94 (for the Hindcast time series) to 0.91 (for the 4-day forecast time series). This small decrease in correlation indicates that the initial state has a small effect. In the second step the lag correlation between the Hindcast (H time series) and the Forecast (F0 to F4 time series) is computed independently for each system. In this case the mean correlation decreases from 0.98 (correlation between H and F0 time series with 1-day lag) to 0.83 (correlation between H and F4 with 5-day lag). Even though the correlation is still high, this stronger decrease indicates that atmospheric forcing has a greater effect in comparison with the initial state. A second diagnostic is based on the error growth computed with the standard deviation of the forecast error, normalized with the standard deviation of the observations (Fig. 11). For the atmospheric variables, the main error is displayed by the fresh water flux which does not drive the variability of the mixed-layer depth in our case, as mentioned before. The normalized standard deviation becomes greater than 1, signifying that for the 1-day forecast

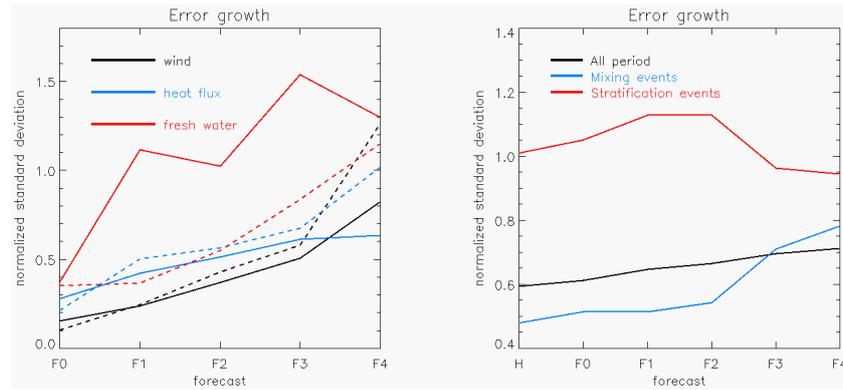


Figure 11. Standard deviation of the forecast error normalized by the standard deviation of the observations. Top panel: atmospheric fields (wind in black, heat flux in blue and fresh water in red) where analyses are considered as observations. The solid line is for May 2013 and the dashed line considers only the mixing events (M1–M3). Bottom panel: ocean mixed-layer depth forecast (for all the systems), in black for May 2013, in blue only during the mixing events (M1–M3) and in red during the stratification events (S1–S3).

the error variance is greater than the observation variance. An unexpected decrease of the error happens for F2. First, current numerical weather prediction systems have difficulties to produce realistic water fluxes over the ocean, in analysis mode as well as in forecast mode. Second, water fluxes may vary a lot inside a given day, or between two instances of a weather forecast. In consequence, errors in water fluxes averaged over 1 day may behave this way due to pure random effects, and a bigger sample may be necessary in order to derive robust statistics for this variable. For the wind field, which in this case is the more important, this ratio is smaller in comparison with the other forcing fields (heat and fresh water fluxes). The difference between the forecast over the entire month and that only over the mixing events (illustrated by the dashed line on the top panel in Fig. 11) is small except for the 4-day forecast. For the mixed-layer forecast (bottom panel in Fig. 11) considering the entire period there is a small linear increase in the normalized standard deviation which generally remains less than 1 even for the 4-day forecast. The link with the error growth for the wind fields can be made by considering that the largest increase in the error for the 4-day forecast will have an effect on the longer-length forecast of the mixed layer (typically for the 5 or 6 days which are not included in this study). Taking only the mixing events into account, the normalized standard deviation is stable for the first 3 days and then increases. It should be noted that during the stratification events the normalized standard deviation for the mixed layer is greater than 1. This is explained by the fact that, in a stratified ocean the error and the mixed-layer depth have the same amplitude and a very small variation in the mixed layer gives rise to a large effect for this ratio. As we see in Figs. 9 and 10, there is also a strong spatial variability in the mixed layer which is not driven by atmospheric forcing, especially at small scales. Computing the spatial standard deviation in the small $1/2^\circ \times 1/2^\circ$ box for all the systems independently, we show that uncertainty

at this small scale is as great, or even greater, than the uncertainty estimated as the standard deviation of all systems and all forecast lengths spatially filtered in the $1/2^\circ \times 1/2^\circ$ box. This standard deviation can reach 50 to 60 m during the month but the available observations are insufficient to quantify this variability in small spatial scale. To understand the initial state differences, an experiment without data assimilation (Atl12 free) was performed and assimilation statistics between systems were compared. The Atl12 free experiment, driven by the best atmospheric forcing, simulates the mixing and stratification events (not shown); the timing of these events is in good agreement with the Atl12 simulation but the amplitude is quite different. The M1 event is too deep and S1 insufficiently stratified, the S2 stratification occurs more quickly and the M3 mixing is insufficiently deep. Statistical results are shown in Fig. 4 where we see that the correlation is still high (0.86), of the same order of magnitude as the 1-day forecast, and the RMSE is comparable with the 2-day forecast. However, the standard deviation is greater than all the Atl12 estimates, showing that data assimilation has a significant effect on the initial state and particularly the stratification which conditions the intensity of the mixing or stratification forecast. Figure 12

shows the SLA increments computed for the three systems (note that there is no data assimilation in the Ibi36 system, which is not presented here). Our area of interest (48.5° N and 16.2° W) is along a well marked front present in all analyses. Positive increments in the northern part and negative in the southern part are deduced from the analysis at $1/4$ and $1/12^\circ$ even though the spatial scales are different with increments containing more mesoscale features at $1/12^\circ$. This front is more intense in the $1/12^\circ$ solution and is further north in Glo12 by comparison with Atl12. These centimetre-scale differences affect the circulation and especially the circulation around mesoscale structures as can be seen in the daily mean for 13 and 16 May (Figs. 9 and 10). The tempera-

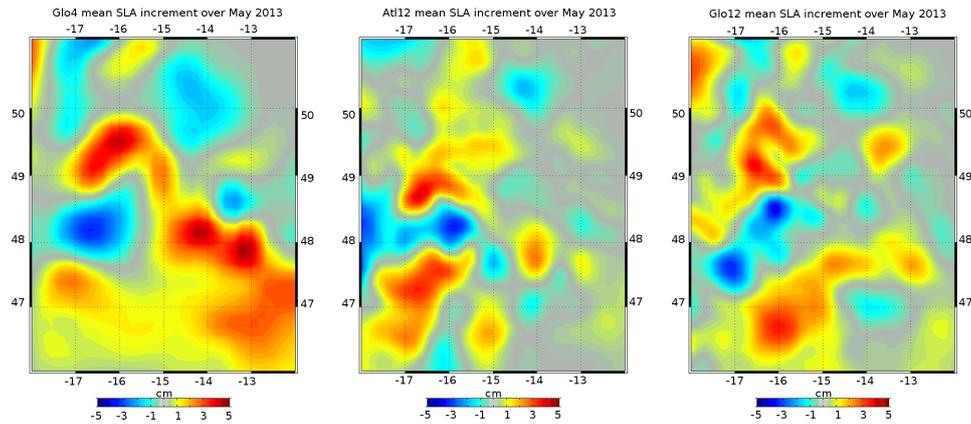


Figure 12. Mean SLA increment (in cm) computed over May 2013 for Glo4 (left), Atl12 (middle) and Glo12 (right) systems.

ture increments presented in Fig. 13 illustrate the correction computed on the temperature profiles during May 2013 as a result of the data assimilation method. Differences between the three systems are noticeable. As already mentioned, the Atl12 system is the closest to observations with a positive increment around $0.1\text{ }^{\circ}\text{C}$ at the surface and a negative increment of the same order of magnitude at 150 m depth. This correction tends to stratify the ocean (warming in the surface layer and cooling at the base of the mixed layer), as is expected given the previous results (Fig. 7). For the Glo12 system, the temperature increment is negative from the surface down to 150 m depth, but also with greater cooling at the base of the mixed layer than in the surface layer. The effect can be considered equivalent to that for the Atl12 system, neglecting the bias. In Glo4, the increment profile is quite different: in the top first 30 m there is a cooling of the mixed layer and then increments re-stratify the ocean from 30 to 150 m just as in the other systems. The dashed lines in Fig. 13 give the envelope of the five increments available over the month of May, computed as plus or minus one standard deviation departure from the mean (we recall that the analysis cycle is one week and in this case we use the five analyses using observations for May 2013). This illustrates the large variability in this increment during this month. This might be expected because of the rapid strong mixing and re-stratification events observed during this month. The conclusion of this part is that evidence of the link between the wind and the mixed-layer forecast is clearer than for the initial state in a complex and non-linear operational system. It is difficult to discriminate the two sources of error using 5-day forecasts that happen to be correlated with the initial state. However, using the Atl12 free simulation which is sufficiently far from initialization with data assimilation (made in March), we were able to show that the effect of data assimilation on the initial state including mesoscale processes and ocean stratification is actually significant. Model physics (vertical mixing scheme) and resolution (from $1/4$ to $1/36^{\circ}$) also play a crucial role; they

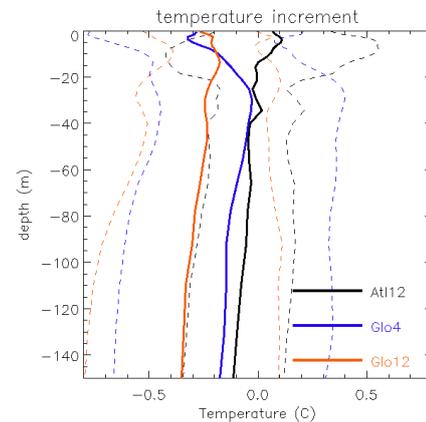


Figure 13. Mean temperature increment (solid line) and \pm one standard deviation departure from the mean (dashed line) for May for the three systems (Glo4 in blue, Atl12 in black and Glo12 in red).

have been discussed and their effects quantified in terms of the statistics generated by the operational systems available.

5 Conclusions

This study focuses on a small area in the Northeast Atlantic during May 2013. Several conditions are met to obtain robust results:

1. a large number of temperature profiles (74) in a small area with a high sampling frequency over the month (more than one per day);
2. available daily forecasts with four operational ocean forecasting systems containing differences in horizontal resolution from $1/4$ to $1/36^{\circ}$, initialization method, vertical mixing scheme, atmospheric forcing, and so on;
3. a strong variability in the mixed-layer depth during the month with alternating mixing and stratification events;

4. a strong link between atmospheric forcing and ocean response.

As a result of all these conditions, we have shown how operational oceanic systems can provide a mixed-layer forecast, and we have quantified the quality of these forecasts with commonly used diagnostics. The mean bias of the mixed-layer depth forecast over the month is around a few metres (usually less than 5 m) and is quite stable with the forecast length; the mixed-layer depth RMSE increases with the forecast length but remains less than 20 m. The accuracy of the mixed-layer depth depends on the vertical grid of the models, and in Mercator Ocean forecasting systems near 100 m, two levels are more than 10 m apart. Improving the vertical resolution of the models could significantly improve the accuracy of the mixed-layer depth estimates. The temporal correlation between observation and forecast is usually greater than 0.85 and slowly decreases with forecast length. The skill score shows the benefit of comparing the forecast with the persistence. These statistics are also useful in comparing the performance of the systems from the best to the worst in terms of forecast ability. In our case we have shown that Glo4, which is the system with the lowest resolution, gives the worst results and Ibi36, which has the highest resolution, gives the best results closely followed by the Atl12 system. This paper concentrates on temporal variability since, with the observations available, it is not possible to estimate a spatial distribution of the mixed-layer depth. We have shown that temporal variability is mainly driven by atmospheric forcing (especially the wind field) and that the model forecast is often close to the observations with good agreement of the temporal sequence of the mixing and stratification events in the observations and forecasts. Note that a ~ 2 -day lag between a strong wind event and the maximum of mixed-layer depth is observed, and consequently missing this event on the first day of the wind forecast generates an error in the mixed-layer depth forecast.

The availability of four systems providing daily forecasts gives the opportunity to build an ensemble forecast associated with an estimate of the uncertainty of the mixed-layer depth. These systems have been developed by Mercator Ocean under the MyOcean project, the ocean part of the European Copernicus programme, and have been operated in real time since the end of April 2013 (V3 of MyOcean service). Other ocean forecast products could also have been used to increase the number of members in the ensemble, but for this study we chose to use only these four forecasts to separate the effects of atmospheric forcing and initial state. First results show the benefit of the mean or the median of the members as forecast. In our case this ensemble estimate is close to the best forecast, and sometimes this estimate is the best (for example the best correlation for the 1-day forecast is obtained with the median state and with the mean for the 4-day forecast). Computing the same statistics, removing each individual forecast one by one, is a good way to estimate

each contribution in the ensemble. We have shown that after removing the worst forecast, which systematically degraded the mixed-layer depth estimation, the mean is always better than each individual forecast for every forecast length. Using other operational forecasts, it will be now useful to introduce into the ensemble ocean estimates computed with other atmospheric forecasts, as for example, the product available in MyOcean provided by the UK's Met Office covering the Northwest shelf (O'Dea et al., 2012), or other global high resolution forecasts such as that provided by Naval Research Laboratory (NRL) (Cummings, 2005). Uncertainty estimates in the mixed layer in this area based on our 4-forecasting systems and 4-day forecast length can reach 50 m during this particular month. The spatial uncertainty for the model in such a small area has the same order of amplitude (~ 50 m). Using the available data an uncertainty of 50 m was also estimated on several dates, though the number of observations might be insufficient to compute a robust level of uncertainty. We have also shown that there is a direct link between the atmospheric uncertainty (especially the wind field) and the mixed-layer depth. The other atmospheric fluxes (net heat and water fluxes) are intrinsically different for each model as they are computed from bulk formulae, which gives more dispersion between all estimates.

Finally we have shown that the temporal variability in the mixed-layer depth when changing from the mixing to the stratification phase is driven by the atmospheric forcing, but the small and meso-ocean scales also have a great local impact. At this smaller scale, resolution, parameterization and assimilation play a role and can impact the forecast score, error or uncertainty. Unfortunately, based on observations the mixing along fronts and around eddies remains difficult to validate properly. The coverage of the in situ observations and the resolution of satellite observations are not sufficient even though the recovery of vertical velocity based on satellite observations is promising (Buongiorno Nardelli et al., 2012) and though observations of water colour provide high-resolution estimates of ocean parameters directly affected by the vertical mixing. However, the effect of horizontal circulation, particularly around eddies or along strong fronts, is illustrated by Figs. 9 and 10. The Ibi36 model, having the highest horizontal resolution, is able to resolve mesoscale eddies that induce patterns of convergence and mixing that are not present in the coarser horizontal resolution systems. The Ibi36 model benefits from model tuning (GLS mixing scheme, explicit tides, higher-resolution atmospheric forcings) that are not yet implemented in the basin scale and global model configurations such as Atl12, Glo12 and Glo4. Further sensitivity studies would be necessary in order to quantify the effect of each individual improvement of Ibi36 with respect to Atl12 or Glo12. Eventually as a test platform for further developments of a high-resolution global system, the Ibi36 forecasting system proves to be successful in reproducing the mixed-layer depth and its response to atmospheric forcing. Future development of the opera-

tional oceanic forecasting systems will be crucial in improving forecasts of oceanic parameters or processes such as the mixed-layer depth. Within the scientific community, work is in progress to include data assimilation of new types of observation (such as ocean colour and, in the near future, SWOT high-resolution sea surface height observations), to increase horizontal and vertical resolution, to improve vertical mixing models and parameterizations, to improve ocean–atmosphere interaction due to coupling and to provide better estimates of the uncertainties based on ensemble techniques. On the short term, Mercator Ocean systems will be improved by using choice already done for Ibi36 as the full resolution of the atmospheric forcing at $1/8^\circ$ in place of a $1/4^\circ$ interpolation, and by modifying the mixing length to 10 m in the TKE mixing scheme or by implementing the GLS mixing scheme (Reffray et al., 2014). A better vertical resolution could also improve the MLD forecast, as well as introducing the mixing due to waves.

Acknowledgements. This research was supported by the MyOcean2 European project and is based on MyOcean products. The authors wish to thank collaborators contributing to the development of the ocean forecasting systems under the MyOcean project, the NEMO consortium and the data centres at CORIOLIS and BODC which disseminate the in situ glider observations collected under the OSMOSIS project.

Edited by: A. Schiller

References

- Balmaseda, M. and Anderson, D.: Impact of initialization strategies and observations on seasonal forecast skill, *Geophys. Res. Lett.*, 36, L01701, doi:10.1029/2008GL035561, 2009.
- Béranger, K., Drillet, Y., Houssais, M. N., Testor, P., Bourdalle-Badie, R., Alhammoud, B., Bozec, A., Mortier, L., Bouruet-Aubertot, P., and Crepon, M.: Impact of the spatial distribution of the atmospheric forcing on water mass formation in the Mediterranean Sea, *J. Geophys. Res.*, 115, C12041, doi:10.1029/2009JC005648, 2010.
- Buongiorno Nardelli, B., Guinehut, S., Pascual, A., Drillet, Y., Ruiz, S., and Mulet, S.: Towards high resolution mapping of 3-D mesoscale dynamics from observations, *Ocean Sci.*, 8, 885–901, doi:10.5194/os-8-885-2012, 2012.
- Cailleau, S., Chanut, J., Lellouche, J.-M., Levier, B., Maraldi, C., Reffray, G., and Sotillo, M. G.: Towards a regional ocean forecasting system for the IBI (Iberia-Biscay-Ireland area): developments and improvements within the ECOOP project framework, *Ocean Sci.*, 8, 143–159, doi:10.5194/os-8-143-2012, 2012.
- Cummings, J. A.: Operational multivariate ocean data assimilation, *Q. J. Roy. Meteorol. Soc.*, 131, 3583–3604, 2005.
- De Boyer Montégut, C., Madec, G., Fischer, A. S., Lazar, A., and Iudicone, D.: Mixed layer depth over the global ocean: An examination of profile data and a profile-based climatology, *J. Geophys. Res.*, 109, C12003, doi:10.1029/2004JC002378, 2004.
- Drillet, Y., Garric, G., Le Vaillant, X., and Benkiran, M.: The dependence of medium range northern Atlantic Ocean predictability on atmospheric forecasts, *Journal of Operational Oceanography*, 2, 43–55, 2009.
- Giordani, H.: Dynamique des couches limites oceanique et atmospherique marine, Habilitation à diriger des recherches, Université Toulouse III-Paul Sabatier, 2011.
- Giordani, H., Caniaux, G., Prieur, L., Paci, A., and Giraud, S.: A 1 year mesoscale simulation of the northeast Atlantic: Mixed layer heat and mass budgets during the POMME experiment, *J. Geophys. Res.*, 110, C07S08, doi:10.1029/2004JC002765, 2005.
- Goni, G. J. and Trinanes, J. A.: Ocean thermal structure monitoring could aid in the intensity forecast of tropical cyclones, *Eos, T. Am. Geophys. Un.*, 84, 573–578, doi:10.1029/2003EO510001, 2003.
- Keerthi, M. G., Lengaigne, M., Vialard, J., de Boyer Montégut, C., and Muraleedharan, P. M.: Interannual variability of the Tropical Indian Ocean mixed layer depth, *Clim. Dynam.*, 40, 743–759, 2013.
- Lavigne, H., D’Ortenzio, F., Migon, C., Claustre, H., Testor, P., Ribera d’Alcalà, M., Lavezza, R., Houpert, L., and Prieur, L.: Enhancing the comprehension of mixed layer depth control on the Mediterranean phytoplankton phenology, *J. Geophys. Res.*, 118, 3416–3430, doi:10.1002/jgrc.20251, 2013.
- Lellouche, J.-M., Le Galloudec, O., Drévilion, M., Régnier, C., Greiner, E., Garric, G., Ferry, N., Desportes, C., Testut, C.-E., Bricaud, C., Bourdallé-Badie, R., Tranchant, B., Benkiran, M., Drillet, Y., Daudin, A., and De Nicola, C.: Evaluation of global monitoring and forecasting systems at Mercator Océan, *Ocean Sci.*, 9, 57–81, doi:10.5194/os-9-57-2013, 2013.
- Lengaigne, M., Menkes, C., Aumont, O., Gorgues, T., Bopp, L., and Madec, J.-M. A. G.: Bio-physical feedbacks on the tropical pacific climate in a coupled general circulation model, *Clim. Dynam.*, 28, 503–516, 2007.
- Lermusiaux, P. F. J., Chiu, C.-S., Gawarkiewicz, G. G., Abbot, P., Robinson, A. R., Miller, R. N., Haley, P. J., Leslie, W. G., Majumdar, S. J., Pang, A., and Lekien, F.: Quantifying Uncertainties in Ocean Predictions, *Oceanography*, 19, 90–103, doi:10.5670/oceanog.2006.93, 2006.
- Lenartz, F., Mourre, B., Barth, A., Beckers, J. M., Vandenbulcke, L., and Rixen, M.: Enhanced ocean temperature forecast skills through 3-d super-ensemble multimodel fusion, *Geophys. Res. Lett.*, 37, L19606, doi:10.1029/2010GL044591, 2010.
- Madec, G. and the NEMO team: NEMO ocean engine, Note du Pôle de modélisation, Institut Pierre-Simon Laplace (IPSL), France, No. 27 ISSN, 1288–1619, 2008.
- Mahadevan, A., D’Asaro, E., Lee, C., and Perry, M. J.: Eddy-driven stratification initiates North Atlantic spring phytoplankton blooms, *Science*, 337, 54–58, 2012.
- Maraldi, C., Chanut, J., Levier, B., Ayoub, N., De Mey, P., Reffray, G., Lyard, F., Cailleau, S., Drévilion, M., Fanjul, E. A., Sotillo, M. G., Marsaleix, P., and the Mercator Research and Development Team: NEMO on the shelf: assessment of the Iberia-Biscay-Ireland configuration, *Ocean Sci.*, 9, 745–771, doi:10.5194/os-9-745-2013, 2013.
- Meehl, G. A., Goddard, L., Boer, G., Burgman, R., Branstator, G., Cassou, C., Corti, S., Danabasoglu, G., Doblus-Reyes, F., Hawkins, E., Karspeck, A., Kimoto, M., Kumar, A., Matei, D., Mignot, J., Msadek, R., Navarra, A., Pohlmann, H., Rienecker,

- M., Rosati, T., Schneider, E., Smith, D., Sutton, R., Teng, H., van Oldenborgh, G. J., Vecchi, G., and Yeager, S.: Decadal Climate Prediction: An Update from the Trenches, *B. Am. Meteorol. Soc.*, 95, 243–267, doi:10.1175/BAMS-D-12-00241.1, 2014.
- Morel, A., Huot, Y., Gentili, B., Werdell, P. J., Hooker, S. B., and Franz, B. A.: Examining the consistency of products derived from various ocean color sensors in open ocean (Case 1) waters in the perspective of a multi-sensor approach, *Remote Sens. Environ.*, 111, 69–88, 2007.
- Murphy, A. H.: Skill scores based on the mean square error and their relationships to the correlation coefficient, *Mon. Weather Rev.*, 116, 2417–2424, 1988.
- O’Dea, E. J., Arnold, A. K., Edwards, K. P., Furner, R., Hyder, P., Martin, M. J., Siddorn, J. R., Storkey, D., While, J., Holt, J. T., and Liu, H.: An operational ocean forecast system incorporating NEMO and SST data assimilation for the tidally driven European Northwest shelf, *Journal of Operational Oceanography*, 5, 3–17, 2012.
- Pham, D. T., Verron, J., and Roubaud, M. C.: A singular evolutive extended Kalman filter for data assimilation in oceanography, *J. Mar. Syst.*, 16, 323–340, 1998.
- Pistoia J.: Development of SuperEnsemble Techniques for the Mediterranean ocean Forecasting System, ALMA MATER STUDIUM Università di Bologna, Dottorato di Ricerca in: GEOFISICA Ciclo XXIII Settore Scientifico-Disciplinare di afferenza: GEO/10, 2012.
- Reffray, G., Bourdalle-Badie, R., and Calone, C.: Modelling turbulent vertical mixing sensitivity using a 1-D version of NEMO, *Geosci. Model Dev. Discuss.*, 7, 5249–5293, doi:10.5194/gmdd-7-5249-2014, 2014.
- Ryan, A. G., Régnier, C., Divakaran, P., Spindler, T., Mehra, A., Hernandez, F., Smith, G. C., Liu, Y., and Davidson, F.: GODAE Oceanview Class 4 forecast verification framework: Global ocean inter-comparison, *Journal of Operational Oceanography*, accepted, 2014.
- Scott, R., Ferry, N., Drevillon, M., Barron, C. N., Jourdain, N. C., Lellouche, J.-M., Metzger, E. J., Rio, M.-H., and Smedstad, O. M.: Estimates of surface drifter trajectories in the Equatorial Atlantic: a multi-model ensemble approach, *Ocean Dynam.*, 62, 1091–1109, doi:10.1007/s10236-012-0548-2, 2012.
- Shapiro, G., Chen, F., and Thain, R.: The effect of ocean fronts on acoustic wave propagation in the Celtic Sea, *J. Marine Syst.*, 139, 217–226, doi:10.1016/j.jmarsys.2014.06.007, 2014.
- Tozuka, T. and Cronin, M. F.: Role of mixed layer depth in surface frontogenesis: The Agulhas Return Current front, *Geophys. Res. Lett.*, 41, 2447–2453, 2014.
- Tranchant, B., Testut, C. E., Renault, L., Ferry, N., Birol, F., and Brasseur, P.: Expected impact of the future SMOS and Aquarius Ocean surface salinity missions in the Mercator Ocean operational systems: New perspectives to monitor ocean circulation, *Remote Sens. Environ.*, 112, 1476–1487, 2008.
- Umlauf, L. and Burchard, H.: A generic length-scale equation for geophysical turbulence models, *J. Mar. Res.*, 61, 235–265, 2003.
- Vandenbulcke, L., Beckers, J. M., Lenartz, F., Barth, A., Poulain, P. M., Aidonidis, M., Meyrat J., Arduin, F., Tonani, M., Fraianni, C., Torrisi, L., Pallela, D., Chiggiato, J., Tudor, M., Book, J., Martin, P., Peggion, G., and Rixen, M.: Super ensemble techniques: Application to surface drift prediction, *Prog. Oceanogr.*, 82, 149–167, doi:10.1016/j.pocean.2009.06.002, 2009.
- Xue, Y., Balmaseda, M. A., Boyer, T., Ferry, N., Good, S., Ishikawa, I., Kumar, A., Rienecker, M., Rosati, A., and Yin, Y.: A Comparative Analysis of Upper-Ocean Heat Content Variability from an Ensemble of Operational Ocean Re-analyses, *J. Climate*, 25, 6905–6929, 2012.
- Zhu, J., Huang, B., and Balmaseda, M. A.: An ensemble estimation of the variability of upper-ocean heat content over the tropical Atlantic Ocean with multi-ocean reanalysis products, *Clim. Dynam.*, 39, p. 1001, 2012.