
DEEP LEARNING WEATHER MODELS FOR SUBREGIONAL OCEAN FORECASTING: A CASE STUDY ON THE CANARY CURRENT UPWELLING SYSTEM

Giovanny A. Cuervo-Londoño	Javier Sánchez	Ángel Rodríguez-Santana
Oceanografía Física y Geofísica	Centro de Tecnologías de la Imagen	Oceanografía Física y Geofísica
Aplicada (OFYGA), ECOAQUA	(CTIM), IUCES	Aplicada (OFYGA), ECOAQUA
Universidad de Las Palmas de Gran	Universidad de Las Palmas de Gran	Universidad de Las Palmas de Gran
Canaria	Canaria	Canaria
Las Palmas, Spain	Las Palmas, Spain	Las Palmas, Spain
giovanny.cuervo101@alu.ulpgc.es	jsanchez@ulpgc.es	angel.santana@ulpgc.es

ABSTRACT

Oceanographic forecasting impacts various sectors of society by supporting environmental conservation and economic activities. Based on global circulation models, traditional forecasting methods are computationally expensive and slow, limiting their ability to provide rapid forecasts. Recent advances in deep learning offer faster and more accurate predictions, although these data-driven models are often trained with global data from numerical simulations, which may not reflect reality. The emergence of such models presents great potential for improving ocean prediction at a subregional domain. However, their ability to predict fine-scale ocean processes, like mesoscale structures, remains largely unknown. This work aims to adapt a graph neural network initially developed for global weather forecasting to improve subregional ocean prediction, specifically focusing on the Canary Current upwelling system. The model is trained with satellite data and compared to state-of-the-art physical ocean models to assess its performance in capturing ocean dynamics. Our results show that the deep learning model surpasses traditional methods in precision despite some challenges in upwelling areas. It demonstrated superior performance in reducing RMSE errors compared to ConvLSTM and the GLORYS reanalysis, particularly in regions with complex oceanic dynamics such as Cape Ghir, Cape Bojador, and Cape Blanc. The model achieved improvements of up to 26.5% relative to ConvLSTM and error reductions of up to 76% in 5-day forecasts compared to the GLORYS reanalysis at these critical locations, highlighting its enhanced capability to capture spatial variability and improve predictive accuracy in complex areas. These findings suggest the viability of adapting meteorological data-driven models for improving subregional medium-term ocean forecasting. This also demonstrates the superior flexibility of graph neural networks compared to traditional models, as they can be adapted to new prediction tasks even when originally developed for different purposes.

Keywords Sea surface temperature forecasting; Graph neural networks; Canary Current Upwelling System; Data-driven ocean prediction; Operational oceanography

1 Introduction

Oceanographic prediction is crucial for understanding climate change and supporting sectors like maritime transport, fisheries, and natural disaster management (Bell et al., 2009). It relies heavily on accurately forecasting mesoscale processes due to their environmental and economic impacts. These processes give rise to distinct structures, influencing mean currents and transporting key ocean properties (Falkowski et al., 1991). Despite its importance, predicting mesoscale processes remains challenging for operational forecasts (Treguier et al., 2017; Mourre et al., 2018), as evidenced by the difficulty in forecasting the Gulf of Mexico Loop Current eddy during the 2010 Deepwater Horizon oil spill (Adcroft et al., 2010; Liu et al., 2013).

Traditional oceanographic prediction techniques rely on numerical models that solve physics-based equations. While these models have significantly advanced ocean forecasting, the theoretical foundation—quasi-geostrophic (QG) theory—and the numerical models have inherent limitations that hinder their accuracy in certain contexts.

Although the QG theory initially advanced our understanding of mesoscale dynamics and ocean prediction, it struggles to address strong currents, cannot account for bathymetric features, fails to incorporate surface density gradients, and does not model frontal dynamics (Cushman-Roisin et al., 1990). Despite these theoretical limitations, global ocean circulation models remain the primary tools for oceanographic forecasting, although they contribute to the persistent challenges of operational mesoscale prediction.

On the other hand, numerical ocean prediction (NOP) models—such as NEMO (Madec et al., 2024), which forms the core of reanalysis systems like GLORYS (Jean-Michel et al., 2021) and operational forecast systems like PSY4V3R1 (Lellouche et al., 2018)—remain the current standard for short-term deterministic forecasting. Nevertheless, they cannot accurately represent reality due to diverse constraints such as incomplete understanding of subgrid-scale parameterizations, poorly known forcing fields, insufficient knowledge of interactions with other Earth system components, and restricted computational resources (Sommer et al., 2018).

Additionally, incomplete observations with spatiotemporal gaps and the limitations of data assimilation schemes—still in continuous improvement—prevent these models from fully capturing the ocean’s state. Furthermore, these models do not fully utilize extensive historical data and lack optimization for modern hardware, such as GPUs, which further reduces their efficiency. These limitations highlight the need for novel approaches to improve ocean prediction capabilities.

In recent years, short-term machine learning weather prediction (MLWP) models have emerged in global atmospheric forecasting, surpassing the efficiency and accuracy of traditional numerical systems (Bouallègue et al., 2024). Models such as Pangu Weather (Bi et al., 2023), GraphCast (Lam et al., 2023), Aurora (Bodnar et al., 2024), Neural-GCM (Kochkov et al., 2023), Gencast (Price et al., 2024), or AIFS (Lang et al., 2024) have demonstrated the power of this approach, significantly reducing inference times and computational costs. By leveraging historical data and focusing on spatiotemporal patterns, MLWP models bypass the constraints of incomplete physical understanding and adapt easily to new prediction tasks without architectural modifications (Dueben and Bauer, 2018; Scher, 2018; Bi et al., 2023).

Nevertheless, MLWP models heavily depend on data quality and availability, facing challenges with heterogeneous, sparse, spatially discontinuous, or noisy datasets, which can limit performance in certain contexts. Additionally, the physical consistency of these models deteriorates over long timescales due to spectral bias, which can result in numerical instabilities or unrealistic hallucinations (Chattopadhyay et al., 2023).

While MLWP models perform well in atmospheric systems, their application to oceanography presents distinct challenges. Unlike atmospheric data, which is abundant, spatially continuous, and less influenced by external factors, oceanographic data is constrained by atmospheric interference, which affects observational accuracy and consistency, spatial discontinuities, predominantly induced by the presence of continental landmasses, which significantly disrupt data continuity and turn the training of these models complex.

Recent advancements in machine learning ocean prediction (MLOP) have led to the development of innovative models at global and regional scales. Xihe (Wang et al., 2024), for instance, has been introduced for global ocean forecasting, while SeaCast (Holmberg et al., 2024) and OceanNet (Chattopadhyay et al., 2024) focus on regional applications. Similarly, progress in mesoscale ocean forecasting, including studies on eddy shedding predictability in the Gulf of Mexico (Zeng et al., 2015) and data-driven turbulence forecasting using autoregressive techniques (Chattopadhyay et al., 2021), has shown promising results.

The Canary current upwelling system (CCUS), the only eastern boundary upwelling system with islands, presents a unique prediction challenge. In this region, intense mesoscale activity results from the interplay between oceanic features, coastal regimes, and bathymetry, driving strong mesoscale stirring and creating a highly dynamic environment (Sangrà et al., 2009; Arístegui et al., 1994; Barton et al., 1998). These dynamics create a complex environment for sub-regional ocean prediction. State-of-the-art MLOP models have not been applied to medium-term forecasts in this region, highlighting the need for novel approaches.

This work adapts the GraphCast model (Lam et al., 2023) for sub-regional ocean forecasting to evaluate its performance under ocean-specific conditions. Using sea surface temperature (SST) as a case study, we introduce modifications to handle challenges like spatial discontinuities and satellite-based data, including spatially masked loss functions. We assess whether the model can capture mesoscale features such as upwelling, potentially addressing the limitations of traditional ocean models in resolving frontal dynamics and subgrid-scale processes.

In this context, SST is useful for studying the dynamics of mesoscale structures (Hausmann and Czaja, 2012), as eddies, fronts, and filaments create distinct thermal signatures, i.e., warm/cold rotating cores. This study uses the Copernicus Marine Service L4 SST reprocessed product (1982–2020), spanning 39 years, for a subdomain within the IBI region (Iberia, Biscay, Ireland) to train and validate ocean models. The dataset provides daily, gap-free sea surface temperature fields, derived from multiple intercalibrated satellite sources.

In the experimental results, we compare the ocean-adapted GraphCast-based graph neural network (GNN) and ConvLSTM-based model (Shi et al., 2015) against NEMO-based reanalysis and forecast products (Jean-Michel et al., 2021; Lellouche et al., 2018)—specifically GLORYS and PSY4V3R1—in terms of predictive skill, computational efficiency, and suitability for short- to medium-term SST forecasting. We focus on the complex CCUS region, which features strong mesoscale turbulence, persistent upwelling fronts, and energetic eddies—challenges for traditional quasi-geostrophic models. We evaluate how well ML-based models capture complex dynamics over different forecast times and their sensitivity to mesoscale features in challenging areas like oceanic islands and coastal capes.

Additionally, the study assesses the ability of these models to capture seasonal and interannual variability in coastal upwelling systems, where mesoscale processes often dominate. Model sensitivity to observational errors and initial conditions is also analyzed, highlighting their influence on forecast degradation. Finally, the work explores operational implications and proposes strategies to bridge the gap between data-driven approaches and traditional numerical paradigms, emphasizing integration of physical constraints and improving resolution in coastal environments.

Our model’s superior performance in SST forecasting is evidenced by $>76\%$ and 48% reductions in RMSE over GLORYS at 5- and 10-day lead times, respectively, and by a computational speed nearly $\sim 100\times$ faster (20-day forecasts in 2.3 minutes vs. 7-day forecasts in 4 hours). However, its sensitivity to initial condition errors—reflected in an 8.3% increase in days exceeding instrumental error thresholds (2017–2020)—and the presence of triangular artifacts that raise RMSE variability by $15\text{--}20\%$ at larger scales point to key architectural limitations. These findings underscore the need for hybrid designs that combine GraphCast’s spatial precision with ConvLSTM-like stabilization to achieve a better trade-off between accuracy and operational robustness.

Section 2 presents the study area of the CCUS and describes the L4 SST dataset used in our experiments. Section 3 introduces the predictive approaches we evaluate in this work, including the GLORYS12V1 reanalysis, ConvLSTM architecture, and our proposed graph neural network. The experimental results, presented in Section 4, assess the performance of the models, with particular emphasis on regions exhibiting intense mesoscale activity and the principal capes. Section 5 analyzes the strengths and limitations of the methods, addresses systematic errors in coastal versus open-ocean zones, and evaluates model robustness under dynamic oceanic conditions. Finally, Section 6 synthesizes the key contributions of this work, outlines practical implications for operational oceanography, and proposes future research directions.

2 Area of study and dataset

2.1 The Canary current upwelling region

The study focuses on a subdomain within the Canary Current Large Marine Ecosystem, specifically the Moroccan subregion (Aristegui et al., 2009), which extends from 21°S to 33°N between Cape Sim and Cape Blanc; see Figure 1. This region includes two distinct meridional upwelling zones, as described by Cropper et al. (2014): i) the $21\text{--}26^{\circ}\text{N}$ zone, characterized by strong, permanent upwelling throughout the year; and ii) the $26\text{--}35^{\circ}\text{N}$ zone, where upwelling remains permanent but is weaker, intensifying during summer due to the seasonal migration of trade winds. The contrast between these two upwelling zones, both in intensity and seasonal variability, highlights the complexity of the Moroccan subregion and underscores its importance for accurately predicting sea surface temperature dynamics.

The region experiences upwelling-favorable winds year-round, with peak intensity during summer due to the northward migration of the Azores High (Wooster et al., 1976). Additionally, the region exhibits pronounced mesoscale oceanographic variability driven by geographic heterogeneity, including variations in continental shelf width, prominent capes, and perturbations induced by the Canary Islands, which generate filaments and eddies.

The interplay of interconnected physical processes governs the coastal upwelling dynamics off northwest Africa, particularly around Cape Ghir. Prominent capes, such as Ghir, Sim, and Cantin, serve as critical control points where topography and bathymetry modulate atmospheric and oceanic flows. Wind forcing, the primary driver of upwelling, intensifies Ekman transport between Cape Ghir and Cape Sim, injecting positive relative vorticity (Troupin et al., 2012) and enhancing the upwelling of cold, nutrient-rich subsurface waters through shear-driven turbulent mixing (Estrada-Allis et al., 2023). The Cape Ghir Plateau, a submarine projection of the High Atlas orogeny on land, deflects the coastal jet offshore, inducing a potential vorticity imbalance that drives the formation of the characteristic filament (Hagen et al., 1996).

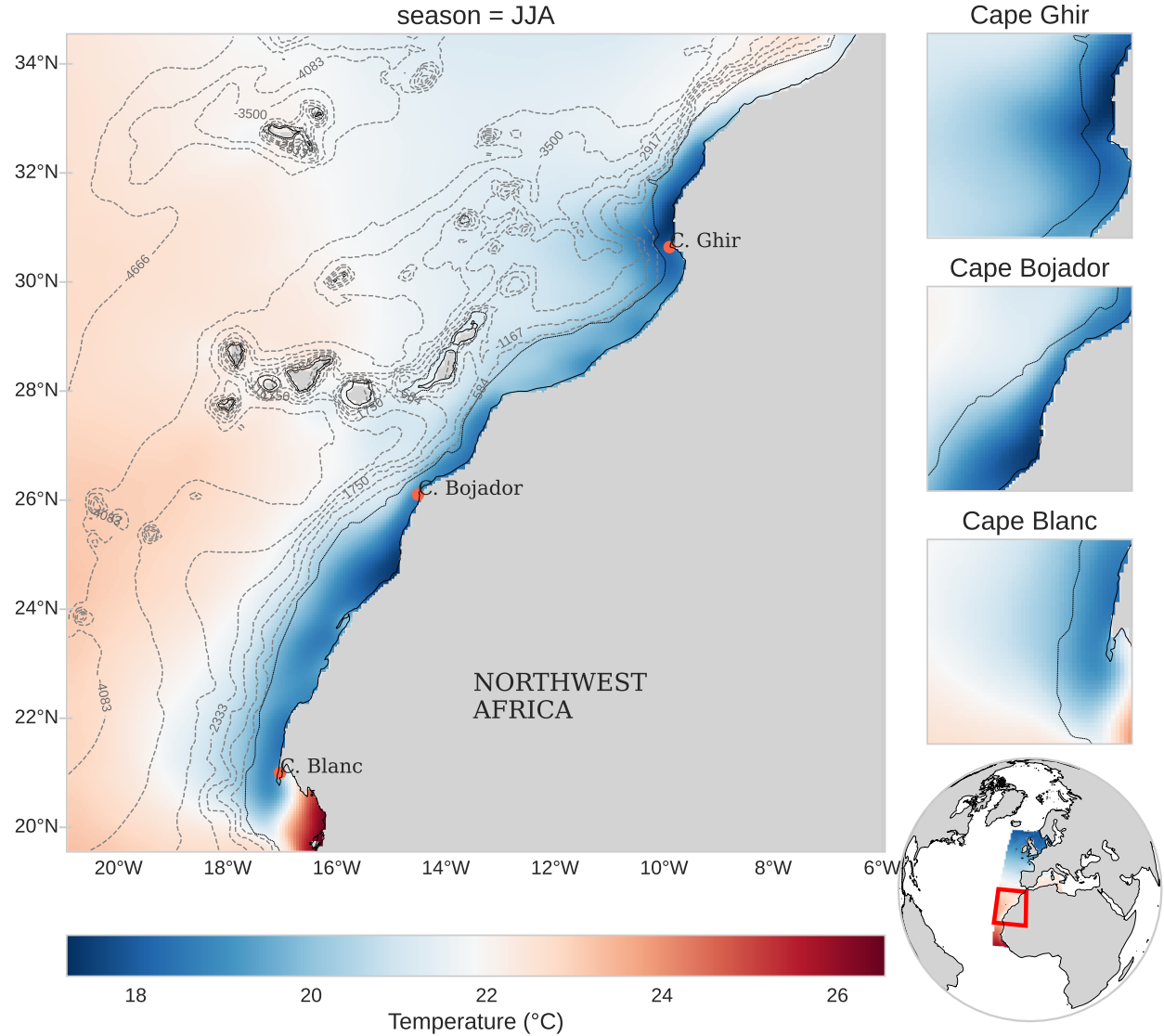


Figure 1: Summer (JJA) climatology of sea surface temperature (°C) in the Northwest Africa region, highlighting the coastal upwelling system. The main panel displays the spatial distribution of temperature along the coast, with three prominent upwelling centers: Cape Ghir ($\sim 30^\circ\text{N}$), Cape Bojador ($\sim 26^\circ\text{N}$), and Cape Blanc ($\sim 21^\circ\text{N}$). Dashed gray contours represent isobaths, indicating bathymetric features relevant to upwelling dynamics. Insets on the right provide zoomed views of each cape to better illustrate localized thermal gradients. The bottom-right globe indicates the IBI (Iberian-Biscay-Ireland) domain of the Copernicus dataset, with a red bounding box marking the subdomain corresponding to the study area shown in the main panel.

This filament displays a dual structure: a cold, surface-intensified core with temperature minima and chlorophyll maxima, surrounded by a broader domain of less dense water influenced by anticyclonic interthermocline eddies (ITEs) (Sangrà et al., 2015). These recurrent ITEs, located north of the filament, strengthen offshore transport through interactions with the upwelling front. Additionally, the irregular bathymetry causes bifurcations in the coastal jet and

forms cyclonic eddies (Hagen et al., 1996), while the subduction of the filament into deeper layers highlights its role in mass and energy export (Sangrà et al., 2015).

Other capes, such as Jubi and Bojador, exhibit similar dynamics where the coastal wind angle and bathymetric irregularities locally intensify upwelling. Atmospheric forcing, topographic constraints, and vorticity adjustments collectively sustain the biological productivity and mesoscale variability unique to the CCUS (Pelegrí et al., 2005).

2.2 High-Resolution L4 sea surface temperature reprocessed

SST data is invaluable for validating models as it evaluates air-sea interactions and vertical mixing while providing insights into the accuracy of model parameterizations and external forcing fields (Mourre et al., 2018). In this study, we use the L4 SST reprocessed product (SST_ATL_SST_L4_REP_OBSERVATIONS_010_026) from the European Union Copernicus Marine Service (CMEMS) for the Atlantic Ocean around Iberia, Biscay, Ireland (IBI), and the northwestern European shelf domain (CMEMS, 2024). It covers nearly 40 years of daily SST data collected by satellites from 1982 to 2020. This high-resolution product, at 0,05 degrees resolution ($\approx 5,55 \text{ km}$), covers the entire IBI domain ($\approx 17\,442\,538 \text{ km}^2$), ranging from $8,93^\circ$ to $61,98^\circ$ latitude and $-20,97^\circ$ to $12,98^\circ$ longitude. It is provided in NetCDF-4 format and is represented in a standard coordinate system (WGS 84/World Mercator). This resource is produced by Ifremer in France and is updated annually.

The L4 product is built from the L3S product SST_ATL_PHY_L3S_MY_010_038 using the inter-calibration method described in Piollé and Autret (2023). Satellite measurements of the SST come from various sources, such as NASA, NOAA, EUMETSAT OSI-SAF, and ESA. These sources are combined using this inter-calibration to create a unified dataset. Each data source includes information about the sensor-specific error statistics to help assess data quality. This information and quality flags are used to identify and select the least reliable data.

For each day, a correction for SST values is estimated to account for discrepancies between the satellites to ensure a consistent daily dataset. For each satellite, a large-scale bias field is determined by comparing the observations of the satellite with the daily reference field. This bias is then smoothed using a Gaussian filter, which helps reduce noise and other irregularities. The smoothed bias is subtracted from the original SST values, resulting in adjusted temperatures that are more accurate and reliable. This adjustment is essential for correcting systemic errors and improving the overall quality of the data.

Once the SST values are adjusted, the single-sensor composite files, L3C, are combined into a multi-sensor composite file, L3S. This merging process ranks sensors based on their accuracy, determined through comparisons with direct measurements. Each cell in the final grid contains data from the best sensor available, ensuring the highest quality SST values.

Our study area is a subdomain within the IBI region, ranging from $19,55^\circ\text{N}$ to $34,525^\circ\text{N}$ and $20,97^\circ\text{W}$ to $5,975^\circ\text{W}$, covering an area of approximately $2\,462\,475 \text{ km}^2$. This region is represented by a grid of 300×300 cells. The temporal range of the data used spans from January 1, 1982, to December 31, 2020, corresponding to a total of 14 245 frames (daily images) and a storage size of 10,25 GB. The final L4 product aims to depict a gap-free daily mean sub-skin SST field in Kelvin at a depth of 20 cm. Therefore, there are no considerable differences between the surface and potential temperatures at this depth.

We preprocessed the dataset by handling missing values for two purposes: first, we used them to calculate a static binary land-water mask, which we smoothed using a Gaussian filter; second, we filled the missing values with the average SST value. We used the average because the model requires complete input data without gaps, and the mean provides a neutral value that avoids introducing strong gradients or artificial patterns.

3 Forecast methods

This section describes the forecasting approaches used to predict SST in the CCUS region, with all models validated against satellite-derived L4 SST data (serving as ground truth). We compare our graph neural network against three baselines: two numerical ocean models and one machine learning approach. First, PSY4V3R1, the operational numerical forecast system that generates daily 10-day predictions using the NEMO platform, represents the state-of-the-art deterministic ocean forecasting. Second, GLORYS, a high-resolution global ocean reanalysis based on the same NEMO framework but enhanced through reanalyzed atmospheric forcing data. Regarding machine learning baselines, we include a ConvLSTM neural network, which combines convolutional layers with long short-term memory (LSTM) (Hochreiter and Schmidhuber, 1997) units to capture spatiotemporal patterns in SST evolution.

Then, we explain our graph neural network in detail, a model that leverages a multiscale graph representation to improve predictions while reducing computational cost. We discuss the modifications to adapt the model for regional

oceanography, including adjustments to the spatial structure and loss function optimization. The specific configurations and hyperparameters for the graph neural network and ConvLSTM models are detailed in Appendix A.

3.1 PSY4V3R1 forecast

Since October 2006, the Mercator Ocean PSY4V3R1 system has provided high-resolution global ocean monitoring and forecasting under CMEMS. With a $1/12^\circ$ ($\sim 9km$) horizontal resolution and 50 vertical levels—offering fine-scale detail in the upper ocean—it captures essential oceanic processes for operational use. Built on NEMO v3.1 (Madec et al., 2008), it integrates high-frequency atmospheric forcing from the European Centre for Medium-Range Weather Forecasts (ECMWF) (Lellouche et al., 2018). Its data assimilation scheme, which combines a reduced-order Kalman filter with 3D-VAR bias correction (Brasseur and Verron, 2006), ingests satellite altimetry, SST, sea ice concentration, and in situ TS profiles (Lellouche et al., 2013).

PSY4V3R1 introduces several key upgrades over its predecessor, PSY4V2. It corrects atmospheric forcing with satellite data, incorporates freshwater runoff from ice sheet melt, and applies a time-varying steric effect to improve sea level representation. The system refines mean dynamic topography with GOCE geoid data (Rio et al., 2011), enhances coastal accuracy with adaptive observational errors, and reduces deep ocean drifts by integrating WOA13v2 climatology.

It also strengthens data reliability through improved TS profile quality control, optimized SSH increments, and assimilation of CMEMS OSI SAF sea ice concentrations (Lellouche et al., 2013). By aligning a 2.2 mm/yr global mass trend with contemporary sea-level rise estimates and deriving background error covariance from bias-corrected simulations, the system improves forecasting stability (Chambers et al., 2017).

The system generates 10-day ocean forecasts, including the runtime day itself, meaning it provides projections for nine days into the future (Galloudec et al., 2024). Updated daily at 00 UTC, it uses the PSY4V3R1 model to predict 3D ocean variables (e.g., temperature, salinity, currents) and 2D variables (e.g., sea level, ice thickness, mixed layer). The atmosphere–ocean coupled NEMO model requires at least 5.5 hours to complete a 6-day forecast when utilizing 864 processors (Thompson et al., 2021). However, in practice, the data assimilation process used in PSY4V3R1 introduces additional computational overhead, so the total runtime for a full cycle may extend to several more hours, even when using a few hundred cores.

3.2 GLORYS12V1 Reanalysis

Jean-Michel et al. (2021) describes GLORYS12 as a high-resolution global ocean reanalysis system based on the ocean and sea ice NEMO models (Madec et al., 2024), starting its simulation in 1991. The system operates NEMO on a quasi-isotropic grid with a $1/12^\circ$ horizontal resolution and 50 vertical levels. The ocean model is coupled with and forced by the ERA-Interim (Dee et al., 2011) atmospheric reanalysis for surface conditions. It also benefits from reanalyzed atmospheric forcing rather than analyses and forecasts, incorporates higher-quality reprocessed observations, and includes refined data assimilation procedures.

GLORYS12 applies the singular evolutive extended Kalman (SEEK) (Brasseur and Verron, 2006) filter method for data assimilation, integrating various sources of information (i.e., satellite sea level anomalies (SLA) (Pujol et al., 2016), satellite SST (Ezraty et al., 2007), and in situ temperature and salinity (T/S) vertical profiles (Cabanès et al., 2013; Szekely et al., 2019)). A 3D-VAR bias correction scheme also estimates large-scale temperature and salinity biases, improving subsurface ocean variability representation.

GLORYS12 utilizes 1 296 processors and completes a 7-day simulation in approximately four hours of computational time. However, the total runtime extends to 14 days because the model employs the incremental analysis update (IAU) method (Lellouche et al., 2013) to assimilate corrections.

3.3 ConvLSTM

The Convolutional LSTM (ConvLSTM) cell with peephole connections, introduced by Shi et al. (2015), is a specialized recurrent neural network (RNN) that combines the capabilities of CNNs to extract spatial correlations with the gating mechanisms of peephole LSTMs to capture temporal dependencies. ConvLSTM is widely used in spatiotemporal prediction tasks, including ENSO forecasting (Mu et al., 2019, 2021), nearshore water level prediction (Yang et al., 2024), and tropical cyclone precipitation nowcasting (Yang et al., 2022).

Researchers initially used RNNs for time-series problems but faced challenges with vanishing and exploding gradients in long sequences (Bengio et al., 1994). LSTM networks introduced gating mechanisms to address these issues, enhancing the handling of extended sequential tasks (Hochreiter and Schmidhuber, 1997; Gers et al., 1999; Gers and

Schmidhuber, 2000). ConvLSTMs further advanced this by incorporating convolutions into their gates, enabling the simultaneous capture of temporal and spatial features.

However, modern LSTM architectures often omit peephole connections, reducing the parameter number without significantly impacting performance (Greff et al., 2017). This adjustment aligns contemporary ConvLSTM implementations more closely with the standard LSTM cells. ConvLSTM replaces traditional matrix multiplications with convolutional operations in input-to-state and state-to-state transitions. This design enables the model to effectively capture spatial features by analyzing local neighboring inputs and states.

The input and forget gates regulate data flow by controlling how new data integrates with previous states to generate an updated cell state, while the output gate determines the current output. Let \mathbf{x}^t be the SST map for a given time instant t , the ConvLSTM equations are given by:

$$\begin{aligned} \mathbf{i}^t &= \sigma(\mathbf{W}_{ii} * \mathbf{x}^t + \mathbf{W}_{hi} * \mathbf{h}^{t-1} + \mathbf{b}_i), \\ \mathbf{f}^t &= \sigma(\mathbf{W}_{if} * \mathbf{x}^t + \mathbf{W}_{hf} * \mathbf{h}^{t-1} + \mathbf{b}_f), \\ \mathbf{g}^t &= \tanh(\mathbf{W}_{ig} * \mathbf{x}^t + \mathbf{W}_{hg} * \mathbf{h}^{t-1} + \mathbf{b}_g), \\ \mathbf{o}^t &= \sigma(\mathbf{W}_{io} * \mathbf{x}^t + \mathbf{W}_{ho} * \mathbf{h}^{t-1} + \mathbf{b}_o), \\ \mathbf{c}^t &= \mathbf{f}^t \odot \mathbf{c}^{t-1} + \mathbf{i}^t \odot \mathbf{g}^t, \\ \mathbf{h}^t &= \mathbf{o}^t \odot \tanh(\mathbf{c}^t), \end{aligned}$$

where \mathbf{i}^t represents the input gate, which determines how much new information from the input sequence \mathbf{x}^t and the previous hidden state \mathbf{h}^{t-1} is allowed into the memory cell. The cell state \mathbf{c}^t acts as a memory that retains relevant information over time and is updated based on the input gate and the forget gate \mathbf{f}^t , which controls the amount of information from the previous cell state \mathbf{c}^{t-1} to retain or discard. The output gate, \mathbf{o}^t , determines how much information from the updated cell state \mathbf{c}^t is passed to the hidden state \mathbf{h}^t , which serves as the final output of the cell at the current time step. The activation functions are the sigmoid, $\sigma(\cdot)$, and the hyperbolic tangent, $\tanh(\cdot)$.

Parameters \mathbf{W}_{ii} , \mathbf{W}_{hi} , \mathbf{W}_{if} , \mathbf{W}_{hf} , \mathbf{W}_{ig} , \mathbf{W}_{hg} , \mathbf{W}_{io} and \mathbf{W}_{ho} represent the weight matrices associated with the input, \mathbf{x}^t , and the hidden state, \mathbf{h}^{t-1} , for each gate, while \mathbf{b}_i , \mathbf{b}_f , \mathbf{b}_g and \mathbf{b}_o are the biases for the respective gates. The Hadamard product (\odot) performs element-wise multiplication, enabling selective gating at each step. While classical LSTMs rely on matrix multiplications, ConvLSTMs replace these operations with convolutions ($*$) within each gate.

3.4 Graph neural network

Our method is based on a GNN model for global medium-range weather forecasting, originally trained on weather reanalysis data, predicting various weather variables globally at a high resolution in under a minute. It is an adaptation of the GraphCast (Lam et al., 2023) model for oceanographic forecasting. This autoregressive model predicts a new state based on two previous time steps. The processing occurs in an underlying multiscale mesh refined from an icosahedron in multiple resolutions. The neural network structure comprises an *encoder*, which embeds the input-grid variables into the mesh nodes, a *processor* that propagates messages through the multiscale mesh, and a *decoder* that maps the forecasting back onto the grid.

The model was originally designed for global atmospheric forecasting and employs two variable types: input and input/target. The model can predict eleven variables: five surface and six atmospheric variables, using data from the ERA5 dataset. These predictions rely on two static input variables, geopotential at the surface and land-sea mask, and five input forcing terms, including solar radiation at the top of the atmosphere and four time-related features.

We adapted this model to predict SST in a local region. This adaptation involved simplifying the model, reducing both the input and input/target variables, and relying on the four time-based features and a single static variable, the land-sea mask. Additionally, to optimize the model for regional use, we replaced the icosahedral mesh with a square curvilinear mesh that represents a curved section of a spherical surface. These modifications highlight the model’s versatility and adaptability for different scales and tasks.

The model uses graphs to simulate the relationship between SST values, represented as discrete cells in a grid. It relies on a bipartite graph $\mathcal{G} = \{\mathcal{V}^g, \mathcal{V}^m, \varepsilon^g, \varepsilon^{g2m}, \varepsilon^{m2g}\}$, made up of two sets of nodes or subgraphs: \mathcal{V}^g , arranged in a grid pattern, and \mathcal{V}^m , structured in a planar and regular mesh. Only the nodes in $v_i^m \in \mathcal{V}^m$ have bidirectional connections via edges $SS_{s,r}^m \in \varepsilon^m$. Additionally, these two sets of nodes are connected by edges, $SS_{s,r}^{g2m} \in \varepsilon^{g2m}$ and $SS_{s,r}^{m2g} \in \varepsilon^{m2g}$, modeling a directional relationship between the grid and the mesh nodes, and vice versa.

This bipartite graph defines local relationships ($\varepsilon^{g2m}, \varepsilon^{m2g}$) between grid node groups $v_i^g \in \mathcal{V}^g$, linked through mesh nodes v_i^m . Additionally, distant multi-scale relationships between these neighborhoods are captured by ε^m connections.

The number of scales, r , is modeled using a multi-mesh configuration $M^r : \mathcal{V}^g$, defined by embedded mesh refinements $\{M^0, M^1, M^2, \dots, M^r\}$. Figure 2 shows an M^2 mesh with red nodes at the coarsest scale with long-range connections, blue nodes at the second scale, and green nodes at the finest scale with short-distant edges.

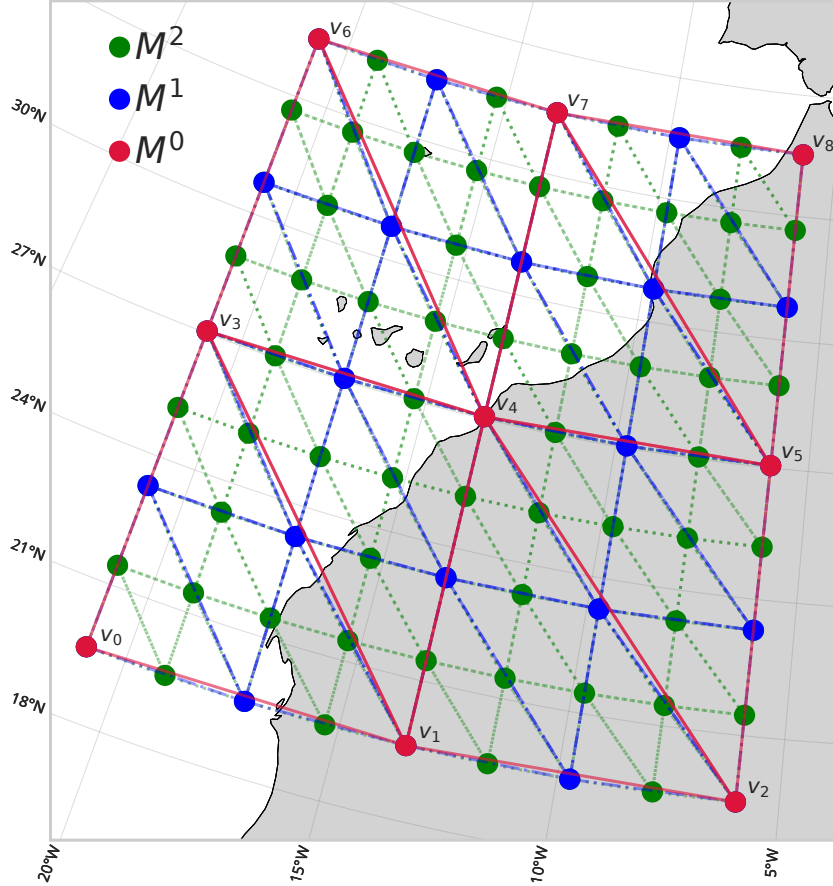


Figure 2: Representation of a multi-mesh with a refinement factor of $r = 2$, consisting of a total of 81 nodes. The nodes are grouped into three resolution levels: 9 nodes belong to M^0 (the coarsest scale, in red), 16 nodes to M^1 (the intermediate level, in blue), and 56 nodes to M^2 (the finest level, in green). We replaced the traditional icosahedral mesh with a curvilinear mesh based on latitude and longitude coordinates. In this approach, each node’s position is explicitly described by its latitudinal and longitudinal values, ensuring a more accurate representation of the Earth’s spherical geometry. New nodes are generated by refining the angular midpoint between existing nodes in spherical coordinates, avoiding distortions from planar approximations and improving spatial accuracy for regional-scale modeling.

The model uses a learnable algorithm called Interaction Network (IN) (Battaglia et al., 2016; Watters et al., 2017) to define how nodes in the graph interact with others. This IN is designed to understand relationships in complex systems (Battaglia et al., 2018; Pfaff et al., 2020; Keisler, 2022). At its core, the IN relies on multilayer perceptrons (MLPs), which in the standard GraphCast implementation typically employ a latent size of 512. However, we use smaller latent sizes since our model forecasts a single oceanographic variable. This reduction in latent dimensionality offers a significant advantage, resulting in a more parsimonious model with fewer parameters. Consequently, the model becomes computationally less demanding, leading to accelerated training and inference.

The IN mechanism facilitates sending messages from one sender node s to another receiver node r , allowing information to flow and update the features of both the nodes and the edges across \mathcal{G} . The *encoder*, *processor*, and *decoder* work through specific message-passing steps within different parts of \mathcal{G} . In the *encoder*, nodes in \mathcal{V}^g send messages to nodes in \mathcal{V}^m , $s^g \rightarrow r^m$. The *decoder* then reverses this process, with nodes in \mathcal{V}^m sending messages back to nodes in \mathcal{V}^g , $r^g \leftarrow s^m$. The *processor* handles messages between nodes in \mathcal{V}^m , $s^m \rightarrow r^m$, performing this step one or more times. Each additional message-passing step increases the number of model parameters, as each step requires an independent IN.

More details about the architecture of the model are given in Appendix B, with the structure of the bipartite graph and the equations of the *encoder*, *processor*, and *decoder*.

4 Results

4.1 Metric-Based Evaluation and Quantitative Analysis of Forecast Skill

The performance of the models was verified against the SST L4 satellite data and assessed through a comprehensive forecast verification framework, detailed in Appendix C, based on four metrics: Root Mean Square Error (RMSE), Anomaly Correlation Coefficient (ACC), Bias, and Relative Activity (RA). These metrics were calculated following the methodology described in Bouall  gue et al. (2024).

Each score was computed independently for each lead time, up to 20 days, using daily temperature fields. The evaluation was performed on a gridded domain of 300×300 points in latitude and longitude, covering the Morocco subregion. For each grid point, the scores were averaged over $365 \times 4 \times 20$ individual forecast realizations, corresponding to daily forecasts generated from 2017 to 2020, the time range of the test set, as shown in Figure 3. This methodology yields robust statistical estimates of model performance for each lead time.

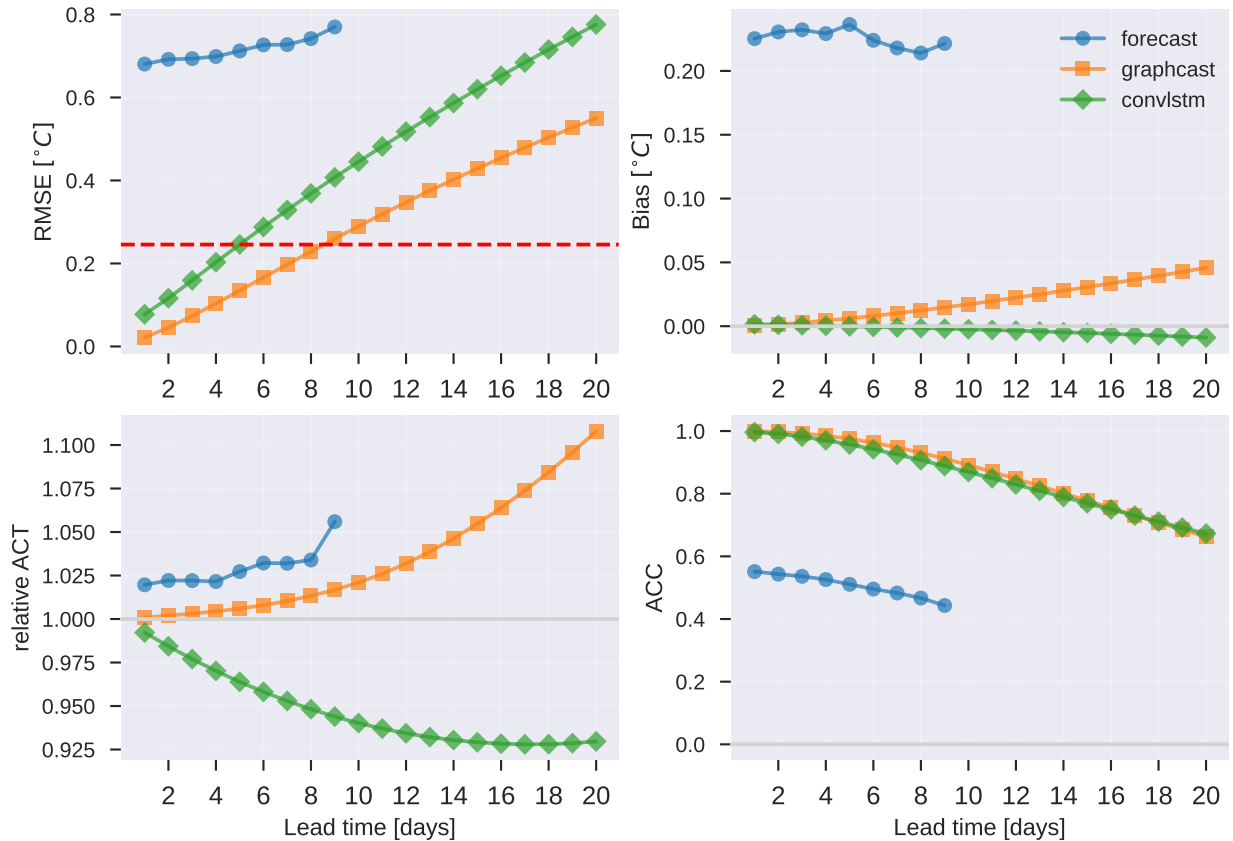


Figure 3: Performance evaluation of the models for SST, verified against satellite L4 observations. The plots display four key metrics over a 20-day forecast period: RMSE, Bias, Relative ACT, and ACC. The models compared are GraphCast (orange), ConvLSTM (green), and PSY4V3R1 (blue), while the red dashed line in the RMSE panel represents the instrumental error of the satellite data ($\pm 0.25^{\circ}\text{C}$).

The SST satellite product used for verification contains an instrumental error for each grid cell within the domain and for each day of the test dataset. The average of this instrumental error is $\pm 0.25^{\circ}\text{C}$ across the entire dataset. This value serves as a reference threshold in the RMSE evaluation, establishing a lower bound for the expected accuracy of the predictive models. In this context, our model surpasses the instrumental error threshold on the eighth day and ConvLSTM on the fifth day, indicating differences in their ability to match the observational accuracy of the satellite product.

A quantitative comparison of the forecast Skill Score (S_{score}) between the two deep learning models was performed using the normalized difference of their scores. This methodology follows the approach described by Geer (2016). The relative RMSE (Eq. 17) and Activity (Eq 18) percentage errors were calculated as:

$$S_{rmse} = \frac{(\overline{\text{RMSE}}_B - \overline{\text{RMSE}}_A)}{\overline{\text{RMSE}}_B} \times 100 \quad (1)$$

where A represents the GraphCast model, and B may represents the reference GLORYS or ConvLSTM models. For the ACC (Eq. 16) and bias (Eq. 21) metrics, which can take negative values, the normalized difference was calculated as:

$$S_{acc} = \frac{(\overline{\text{ACC}}_A - \overline{\text{ACC}}_B)}{(1 - \overline{\text{ACC}}_B)} \times 100 \quad (2)$$

The spatial averages were computed on the 300×300 grid cells in every forecast, and the temporal averages were subsequently calculated for each day. This resulted in 20 values representing the normalized difference for each lead time. Each value represents the average of $300 \times 300 \times 365 \times 4$ data points.

The S_{rmse} and S_{acc} quantify the relative performance of GraphCast compared to ConvLSTM. Those percentages indicate that our model achieves a higher or lower score than ConvLSTM at a given lead time.

Looking at Figure 3, we observe that, for the initial five lead times, GraphCast demonstrates a 81,7% to 38,4% higher S_{acc} relative to ConvLSTM, while for the subsequent ten lead times (6 to 10), the S_{acc} is 30,9% to 7,8% higher. In terms of S_{rmse} , GraphCast exhibits an error reduction between 74,6% and 30,0% in the first ten lead times, and an improvement of 28,4% to 21,7% in the last ten lead times. These results highlight the enhanced predictive skill of the graph model across short and extended forecast horizons.

In terms of RA , GraphCast shows overactivity with a monotonically increasing level relative to observations, while ConvLSTM shows underactivity, with a consistently lower and decreasing level across lead times. Despite these differences, both models maintain a low mean bias across all forecast lead times. ConvLSTM’s bias remains virtually zero ($< 0,009^\circ\text{C}$), whereas GraphCast’s bias increase gradually up to $0,05^\circ\text{C}$ at the last lead time. In both cases, the bias stays well within the instrumental uncertainty. Note that these biases are averages over the spatial domain and all initialization dates yielding a naively results.

4.2 Interannual model performance

We can analyze the precision of the methods using barrier plots, where the y-axis represents the lead time, the x-axis represents the forecast day, and the color scale indicates the RMSE value. These plots depict $365 \times 4 \times 20$ predictions, where each point represents the latitude-weighted spatial average of the 300×300 grid cells for each day. Figure 4 compares the barrier plots of the models from 2017 to 2020.

We can use these barrier plots to assess the percentage of forecast days in which model errors exceeded the mean satellite instrumental error across different lead times and forecast realizations from January 3, 2017, to December 26, 2020. We quantified the proportion of forecast lead time days where RMSE values surpassed the instrumental error threshold (cf. Table 1). Higher percentages indicate periods of reduced model skill, and lower values indicate improved performance.

A notable anomaly occurred on January 1, 2020, during which both models demonstrated unusually low skills. This anomaly coincided with a significant shift in satellite instrumental error, suggesting a potential link between the two events, as illustrated in Figure 4. The magnitude of this instrumental error shift was quantified to range between $3,2\sigma$ to $5,2\sigma$ and was compared against the time series of instrumental error variability. This comparison supports the hypothesis that observational uncertainties contributed to the degraded model performance on that day.

Table 1: Percentage of forecasts exceeding the instrumental error threshold for GraphCast and ConvLSTM (2017–2020). The comparison is based on 28 719 RMSE values, above the instrumental error as the baseline.

Year	N	GraphCast %	ConvLSTM %
2017	7 070	51,0	72,0
2018	7 300	52,0	71,0
2019	7 300	57,0	73,1
2020	7 049	59,3	75,2

Both models recorded the highest percentage in 2020, with 59,3% of RMSE values above the threshold for GraphCast and 75,2% for ConvLSTM. This achieved its lowest value in 2018, at 71%, while GraphCast recorded its lowest in

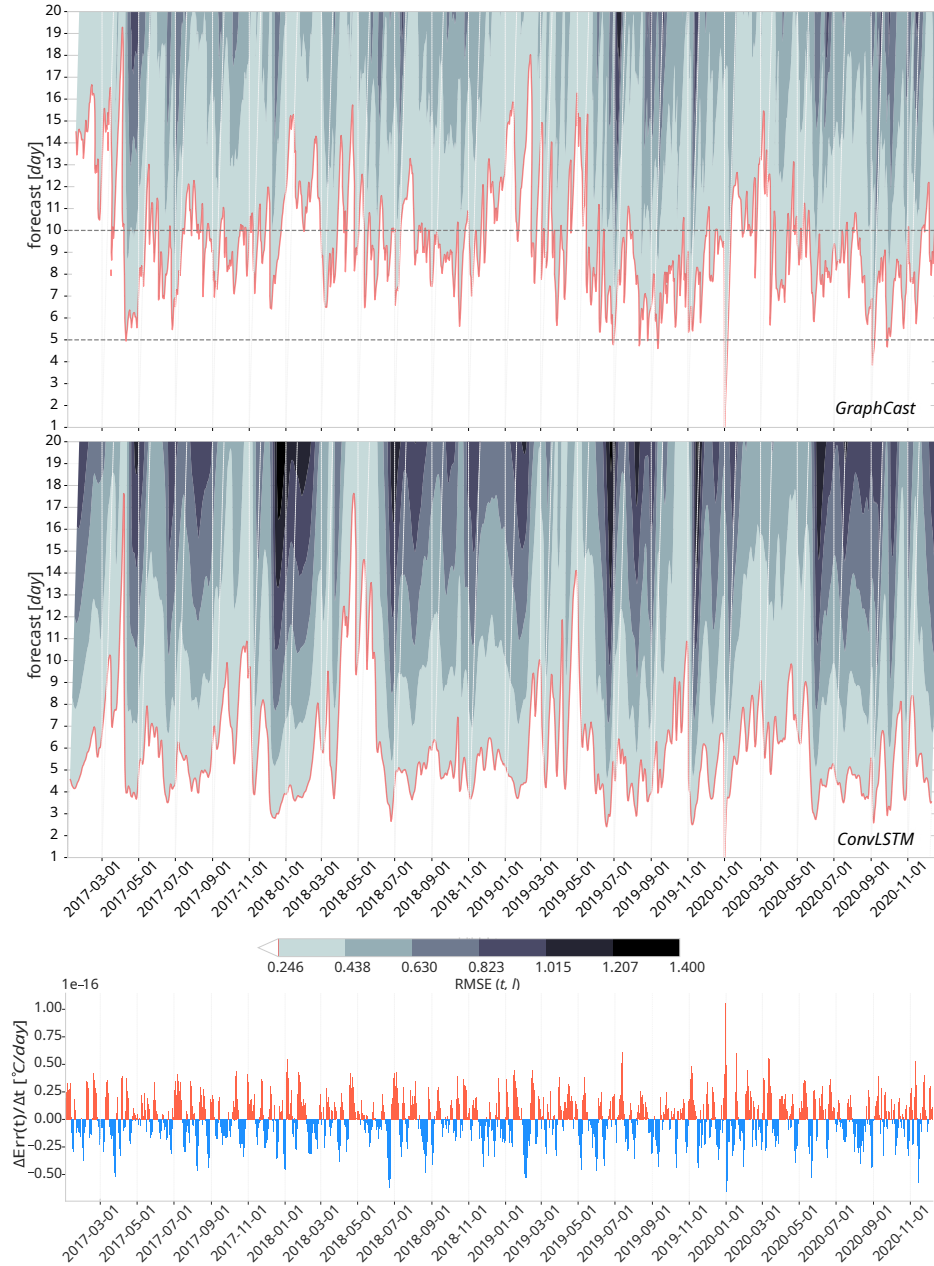


Figure 4: Predictability barrier plots of daily RMSE for SST forecasts at lead times of 1–20 days from 2017 to 2020, comparing GraphCast (top panel) and ConvLSTM (middle panel). Color shading—from pale (low RMSE) to dark (high RMSE)—represents the evolution of forecast error as a function of lead time and forecast initialization date. The solid red line indicates the mean instrumental-error threshold, while dashed grey lines mark 5- and 10-day lead-time references. All RMSE values are validated against L4 satellite SST data. The bottom panel displays time series of daily instrumental-error anomalies ($\Delta \text{Err}(t)/\Delta t$, $^{\circ}\text{C day}^{-1}$), with red (blue) bars indicating increasing (decreasing) errors. Notably, on 1 January 2020, both models exhibit a simultaneous RMSE spike exceeding the mean instrumental-error threshold, coinciding with a sharp positive anomaly in instrumental error—indicating a transient but substantial loss of predictability.

2017 at 51%. This highlights a period of relatively better performance for GraphCast, particularly in 2017, where it was significantly better than ConvLSTM. The results underscored GraphCast’s improved consistency and accuracy over ConvLSTM across the evaluated years.

GraphCast produced a four-year RMSE average of $0,42^{\circ}\text{C}$ with a standard deviation of $0,13^{\circ}\text{C}$. The model showed a slight reduction in 2018 ($0,40^{\circ}\text{C}$) and the highest increment in 2020 ($0,43^{\circ}\text{C}$). In contrast, ConvLSTM yielded higher RMSE values, with a four-year average of $0,53^{\circ}\text{C}$ and a standard deviation of $0,2^{\circ}\text{C}$. The model exhibited its lowest RMSE in 2019 ($0,52^{\circ}\text{C}$) and the highest value in 2018 ($0,54^{\circ}\text{C}$).

Additionally, we calculated the skill RMSE of GraphCast compared to ConvLSTM, following the same methodology described in the previous section. The S_{rmse} shows that GraphCast consistently outperformed ConvLSTM in the test period, with improvements ranging from approximately 19,4% to 26,5%, with the largest improvement in 2018. GraphCast obtained lower RMSE and more values under the instrumental error threshold, consistently indicating superior predictive performance compared to ConvLSTM in all years.

We compare the RMSE time series of GraphCast for 5-day and 10-day forecast lead times with the RMSE of the GLORYS reanalysis data in the test period, as illustrated in Figure 5. The results demonstrate significant improvements in forecast skill: for the 5-day lead time, GraphCast achieves a 75,5% reduction in RMSE compared to the reanalysis; similarly, for the 10-day lead time, GraphCast shows a 47% lower RMSE than the reanalysis. These findings highlight GraphCast’s superior performance in reducing prediction errors across the test period. Additionally, GraphCast generates a 20-day SST forecast in approximately 140 seconds on a Quadro RTX 4000 GPU with 8 GB RAM, showcasing its remarkable computational efficiency. In contrast, GLORYS, which simulates a comprehensive suite of ocean variables—including salinity, velocity components, and others—takes approximately 4 hours to complete a 7-day simulation.

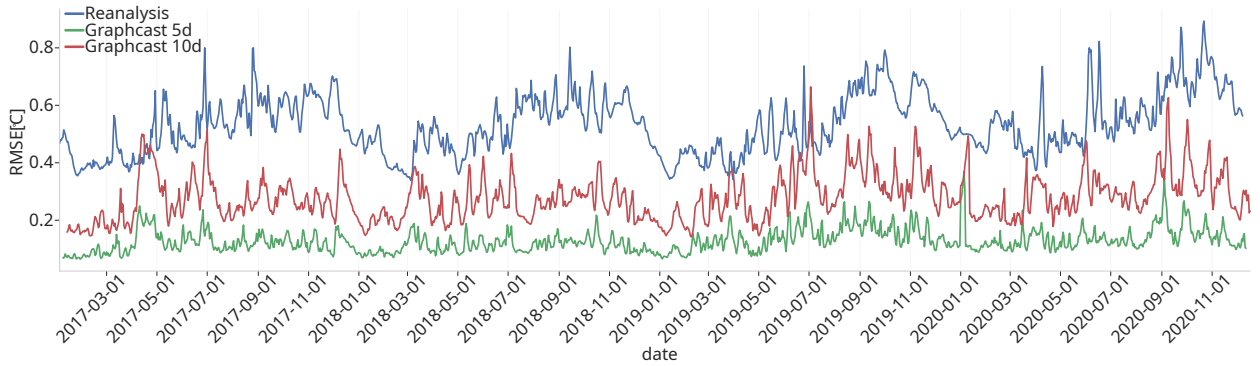


Figure 5: Time series of RMSE for SST predictions over the test period. RMSE values from GraphCast forecasts at 5-day (green) and 10-day (red) lead times are shown alongside those from the GLORYS reanalysis (blue). GraphCast consistently exhibits lower RMSE throughout the period, indicating reduced prediction error relative to the reanalysis.

4.3 Seasonal analysis

To further investigate the seasonal patterns in the model, we used RMSE barrier plots based on a three-month moving window average centered on each day of the year. We apply a 90-day smoothing window to the daily spatially latitude-weighted RMSE values from the barrier plots of each model, spanning the four years of the test set. The window included 45 days before and after the target day, accumulating 90 days, repeated across the four years. Each value represents an average of $300 \times 300 \times 90 \times 4$ RMSE values, resulting in 365×20 data points per day of the year and per lead time. Figure 6 shows the seasonal barrier plots.

The seasonal performance of the models was evaluated by calculating the percentage of days within each season where the RMSE exceeded the mean instrumental error. ConvLSTM recorded the highest rate in DJF and JJA, with 19,1% and 20,2% of RMSE values above the threshold, respectively. GraphCast, on the other hand, showed its peak percentage in JJA and SON, with 15,4% and 15,6% of values above the threshold, respectively. GraphCast demonstrated its lowest rate in DJF, at 12%, while ConvLSTM achieved its best performance in MAM, with 17% of values above the threshold. These results highlight seasonal variations in model performance, with GraphCast showing relatively better consistency across seasons than ConvLSTM.

Additionally, the mean RMSE for each season was computed to assess the seasonal performance of the models (cf. Table 2). GraphCast exhibited the highest RMSE of $0,44^{\circ}\text{C}$ with a standard deviation of $0,1^{\circ}\text{C}$ during SON, while its lowest RMSE of $0,35^{\circ}\text{C}$ with a standard deviation of $0,07^{\circ}\text{C}$ occurred in DJF. In contrast, ConvLSTM produced its highest RMSE of $0,6^{\circ}\text{C}$ with a standard deviation of $0,2^{\circ}\text{C}$ in JJA, and its lowest RMSE of $0,4^{\circ}\text{C}$ with a standard deviation of $0,1^{\circ}\text{C}$ in MAM. These findings align with the previous analysis, reinforcing that GraphCast consistently outperforms ConvLSTM across seasons, particularly in winter (DJF), where it achieves the lowest rate.

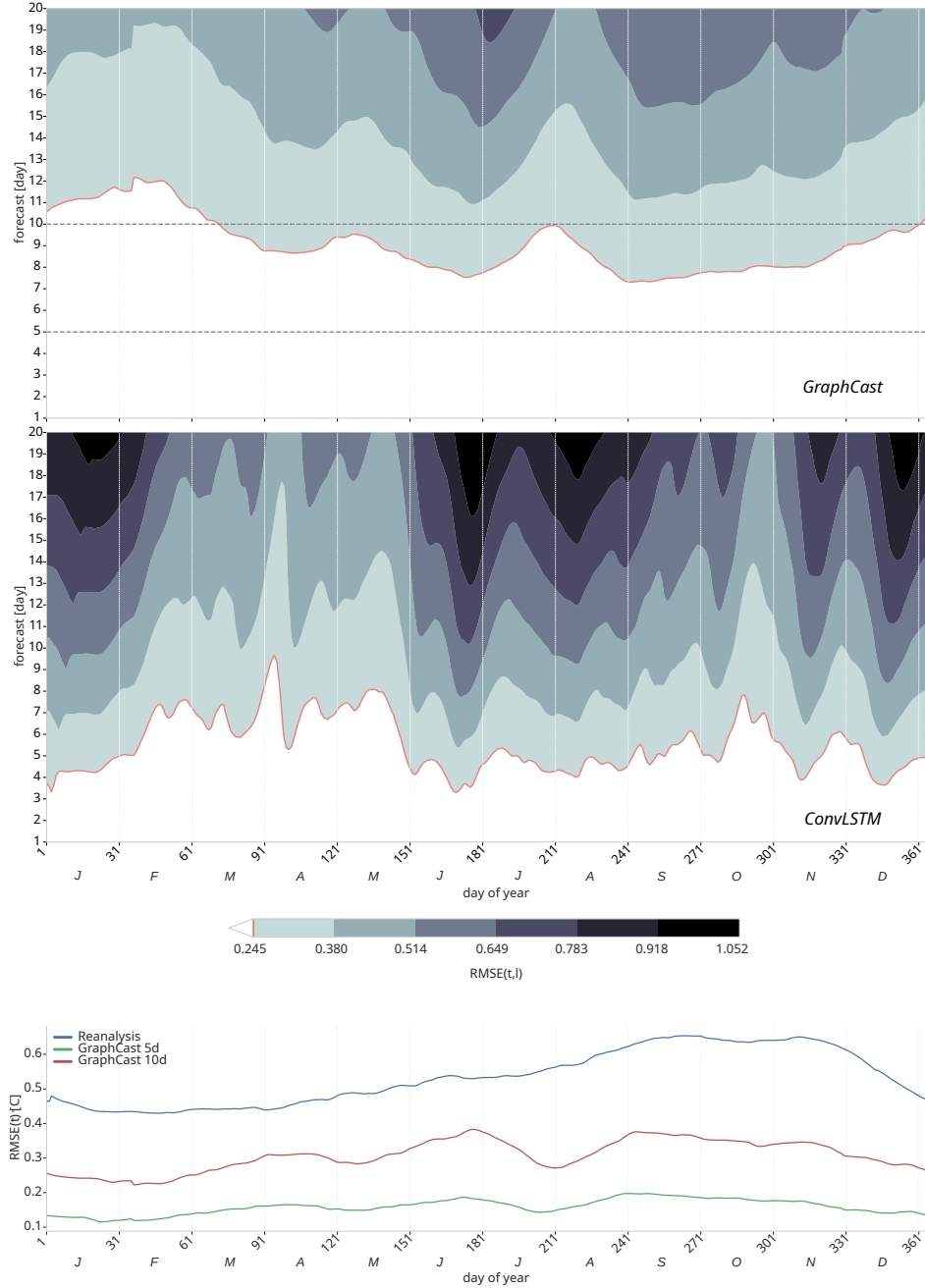


Figure 6: Seasonal predictability barrier plots showing daily RMSE for SST forecasts at lead times of 1–20 days, averaged using a 90-day moving window (± 45 days) to highlight seasonal patterns in forecast skill. The top panel corresponds to GraphCast and the middle panel to ConvLSTM. Color shading—from pale (low RMSE) to dark (high RMSE)—indicates forecast accuracy as a function of lead time and day of year. The solid red line marks the mean instrumental-error threshold, and dashed grey lines indicate 5- and 10-day lead-time references. All RMSE values are validated against Level-4 satellite SST data. The x-axis spans the full calendar year (days 1–365), with each point representing the center of a 90-day moving window, averaged across four years (2017–2020) to capture the seasonal cycle of forecast skill. The bottom panel presents the seasonal cycle of RMSE at selected lead times (5 and 10 days) for GraphCast, along with the corresponding RMSE from the reanalysis (blue), providing a reference for error magnitude across the year. These time series correspond to horizontal slices through the top panel and emphasize periods of relatively higher or lower model skill.

Table 2: Seasonal RMSE for GraphCast and ConvLSTM models. The table presents the mean (μ) and standard deviation (σ) of RMSE values for each season: December-January-February (DJF), March-April-May (MAM), June-July-August (JJA), and September-October-November (SON). GraphCast and ConvLSTM performance is compared across seasons, highlighting variations in prediction accuracy and consistency.

Season	GraphCast		ConvLSTM	
	μ ($^{\circ}C$)	σ ($^{\circ}C$)	μ ($^{\circ}C$)	σ ($^{\circ}C$)
DJF	0,35	0,07	0,58	0,20
MAM	0,38	0,08	0,40	0,10
JJA	0,43	0,11	0,60	0,20
SON	0,44	0,10	0,50	0,15

Next, we calculate the relative improvement of GraphCast to ConvLSTM by estimating the S_{rmse} for each season and lead time. The results demonstrate that GraphCast consistently outperformed ConvLSTM across all seasons. The most significant improvement was observed in DJF, with a 38,5% reduction in RMSE, while the smallest improvement occurred in MAM, with a rate of 7%. GraphCast also achieved lower RMSE values in the remaining seasons, with 28,4% in JJA and 10,6% in SON compared to ConvLSTM. These findings align with seasonal performance trends identified earlier, further emphasizing GraphCast’s superior accuracy and consistency across all seasons, particularly in winter and summer.

Additionally, the same moving window averaging procedure was applied to the RMSE of the GLORYS reanalysis dataset, verified against satellite SST, to establish a reference climatology for the reanalysis. The seasonal RMSE time series for GraphCast at 5-day and 10-day lead times was then compared to the GLORYS reference, providing a comparative view of how well these forecast lead times reproduced the seasonal cycle, as shown in Figure 6c. The normalized difference between GraphCast and GLORYS was calculated to quantify the relative skill of GraphCast in capturing the seasonal cycle compared to the reanalysis data.

We computed the S_{rmse} of GraphCast relative to the GLORYS reanalysis for each forecast lead time to evaluate model performance further. The results demonstrate significant improvements in forecast accuracy: for the 5-day lead time, GraphCast achieved the highest reduction in RMSE during DJF and SON, with a rate of 77,4% for both seasons, while a rate of 73,0% and 75,3% were observed for MAM and JJA, respectively. Similarly, for the 10-day lead time, GraphCast showed the highest skill in DJF and SON, at 51,2% and 50,7%, respectively, with 42,4% and 46,0% for MAM and JJA. This comparison highlights the seasonal variability in model performance and emphasizes the extent to which GraphCast outperforms GLORYS across different seasons and forecast horizons.

4.4 Analysis of the spatial accuracy of models forecasts

In this section, we analyze the accuracy of models in the area of study. We generated RMSE maps for each lead time, where each grid cell represents the temporal average of 365×4 RMSE values, corresponding to the four years of the test set. This averaging procedure results in 20 RMSE maps. Figure 7 shows the corresponding error maps for GraphCast and Figure 8 the maps for the ConvLSTM model.

The first step of the analysis quantifies the percentage of ocean grid cells, out of a total of 49 061, where the RMSE exceeded the mean instrumental error for each lead time and model. This metric provides insight into the spatial extent of regions with significant prediction errors. From the first lead time, ConvLSTM showed 0,9% of cells above the instrumental error threshold, while GraphCast had no cells exceeding this error. By the fourth lead time, GraphCast surpassed the instrumental error in 1,05% of cells, whereas ConvLSTM had a rate of 12,54% of cells above the threshold. Notably, ConvLSTM reached 100% of ocean cells exceeding the instrumental error by the eighth lead time, while GraphCast remained significantly lower at 22,96% for the same lead time. By the twentieth lead time, GraphCast produced 99,01% of ocean cells exceeding the instrumental error. Therefore, GraphCast maintains lower prediction errors across a larger spatial extent and longer lead times than ConvLSTM.

The spatial RMSE average reveals that GraphCast consistently exhibits lower values than ConvLSTM across all lead times. At the first lead time, GraphCast achieves $\mu = 0,02^{\circ}C$ and $\sigma = 0,008^{\circ}C$, while ConvLSTM shows $\mu = 0,05^{\circ}C$ and $\sigma = 0,054^{\circ}C$. However, the standard deviation of GraphCast grows more rapidly than that of ConvLSTM, reaching $\mu = 0,21^{\circ}C$ and $\sigma = 0,094^{\circ}C$ by the eighth lead time, compared to ConvLSTM’s $\mu = 0,36^{\circ}C$ and $\sigma = 0,088^{\circ}C$. Beyond the eighth lead time, GraphCast’s standard deviation remains higher than ConvLSTM, while its mean is lower. At the last lead time, GraphCast achieves a smaller mean ($\mu = 0,5^{\circ}C$) than ConvLSTM ($\mu = 0,76^{\circ}C$), but its standard deviation ($\sigma = 0,218^{\circ}C$) is larger ($\sigma = 0,13^{\circ}C$). Therefore, GraphCast maintains lower prediction errors despite its higher variability in later lead times.

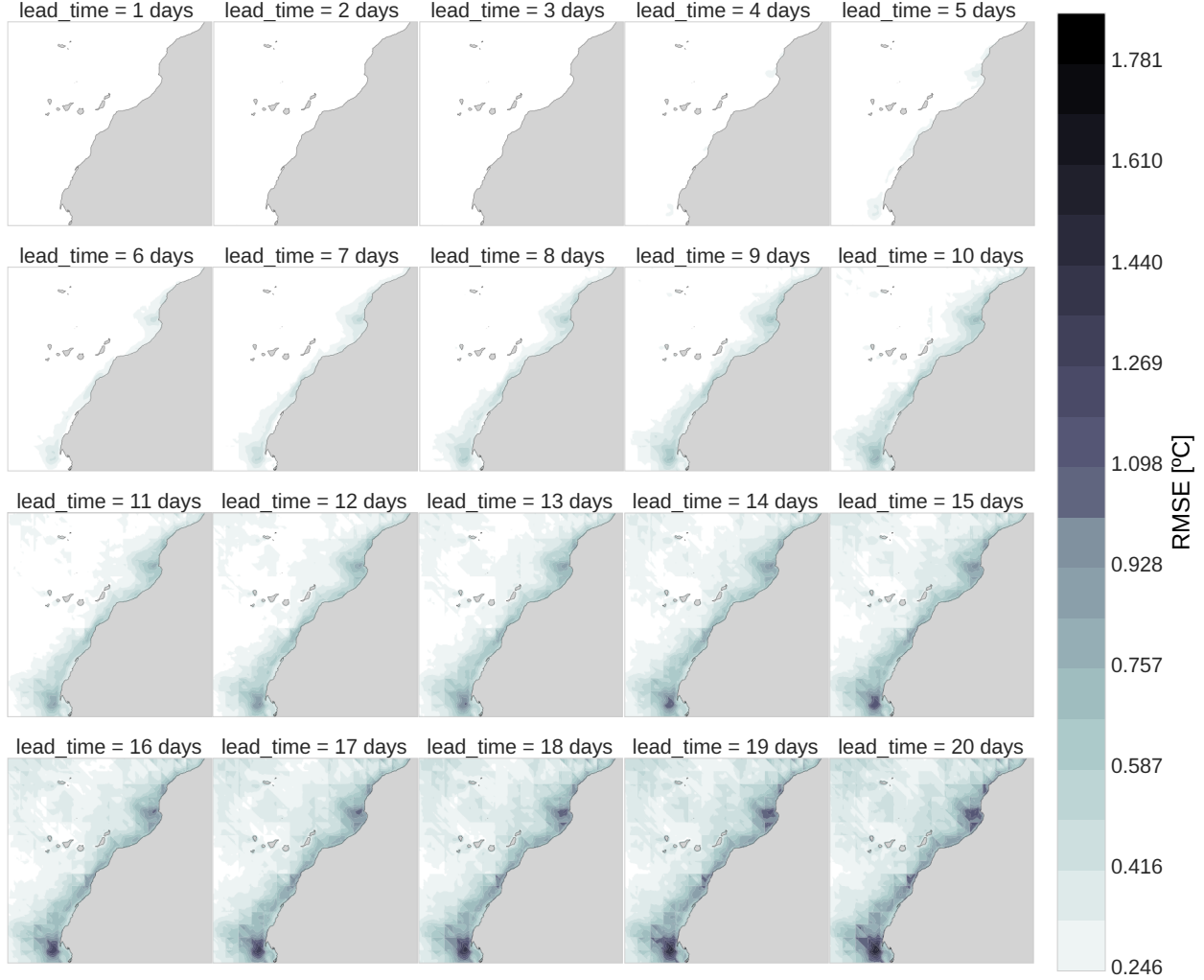


Figure 7: Point-wise RMSE average for 20 lead times forecasts of the GraphCast model. Each map represents the average RMSE at each point of our area of study for 365 days and 4 years of the test set. The images show the results of the 1 to 20 lead times from left to right and top to bottom.

We also computed the spatial average S_{rmse} relative to ConvLSTM. The mean reduction in RMSE across all lead times is 43,86%, with the highest reduction of 62,48% at the first lead time and the lowest reduction of 35,73% at the last lead time. This indicates that GraphCast provides consistently small errors over ConvLSTM across all forecast horizons.

The spatial analysis also allows us to identify regions with high RMSE values. Visual inspection of the maps revealed that the areas with the highest values for both models were concentrated near prominent capes, which collectively accounted for more than 18,6% of the total RMSE across all lead times relative to the entire domain. Specifically, Cape Ghir contributed with a rate of 6,8%, Cape Bojador with 4,0%, and Cape Blanco with 7,8% to the overall RMSE. Despite these challenges, GraphCast demonstrated superior performance in these regions compared to ConvLSTM. On average, GraphCast achieved significant reductions in RMSE relative to ConvLSTM, with improvements of 22,2% at Cape Ghir, 24,9% at Cape Bojador, and 28,5% at Cape Blanco. Therefore, GraphCast can handle complex regional dynamics better than ConvLSTM, particularly in areas with high prediction errors.

Finally, we computed the spatial S_{rmse} relative to the time-average GLORYS reanalysis, with an full domain average RMSE of 0,48 $^{\circ}\text{C}$. This metric quantifies the domain-averaged skill in RMSE between GraphCast and the reanalysis. For the 5-day lead time, the GraphCast average RMSE is smaller by a rate of 74,2% relative to the reanalysis across the entire domain. At the 10-day lead time, the error is reduced by a 44,1%, reflecting diminished but still significant improvements. These results demonstrate GraphCast's ability to outperform the reanalysis benchmark, with higher performance at shorter lead times. The same analysis was applied to each coastal cape and, for the 5-day lead time,

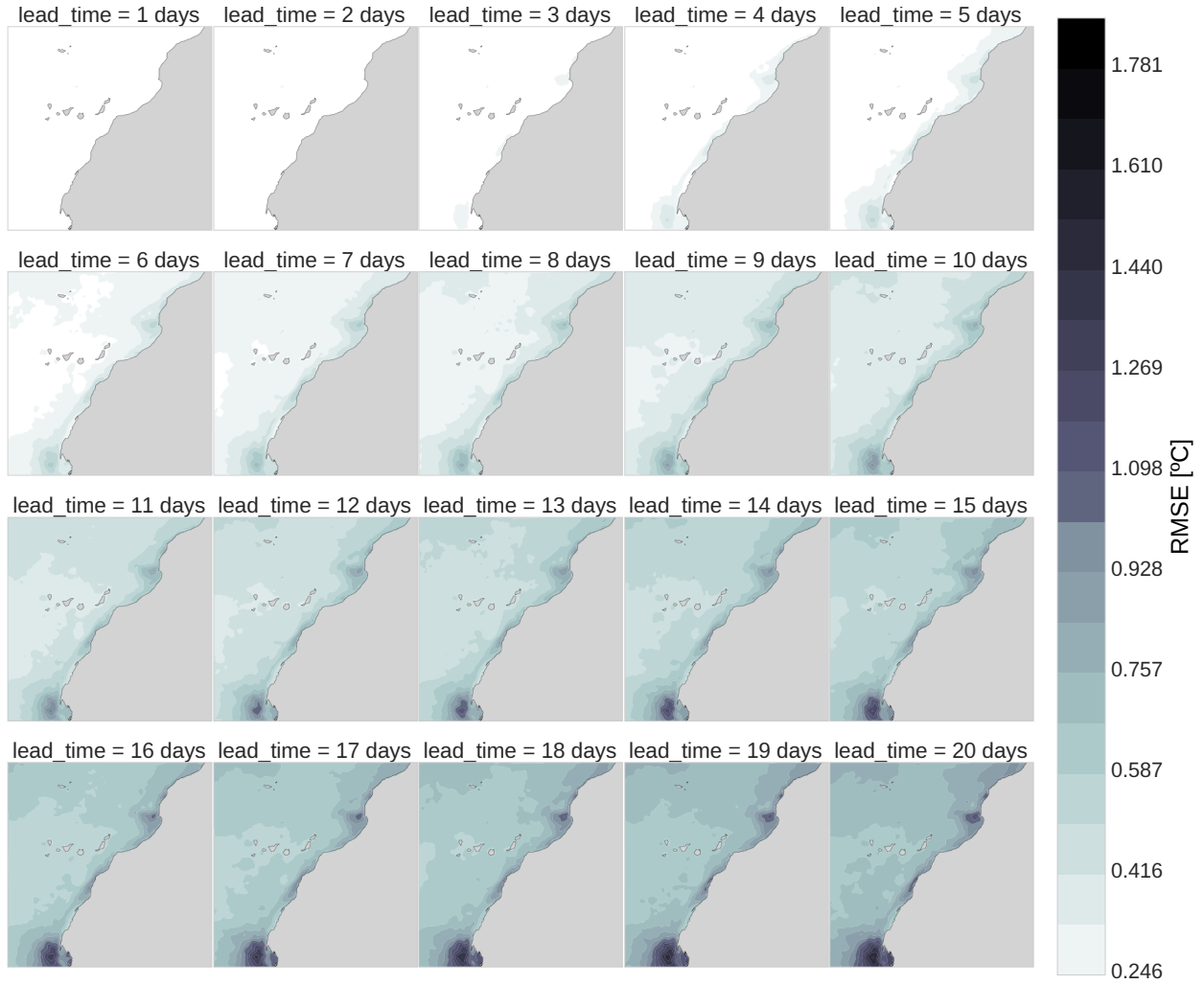


Figure 8: Point-wise RMSE average for 20 lead times forecasts of the ConvLSTM model. Each map represents the average RMSE at each point of our area of study for 365 days and 4 years of the test set. The images show the results of the 1 to 20 lead times from left to right and top to bottom.

GraphCast improved the average RMSE by a 69,7% at Cape Ghir, 78,6% at Cape Bojador, and 78,5% at Cape Blanc relative to the reanalysis. At the 10-day lead time, the reductions were 36,8%, 53,5%, and 51,1%, respectively. This shows that GraphCast provides consistently better estimates than the reanalysis data.

5 Discussion

The results presented in the previous section highlight key differences in performance between the two MLOP models—GraphCast and ConvLSTM—and the GLORYS reanalysis, along with the instrumental error baseline. GraphCast shows consistently lower RMSE and higher ACC values in all lead times, indicating superior predictive capability in the medium-range forecasts.

Between 2017 and 2020, both deep learning models showed an increasing number of days when RMSE surpassed the instrumental error threshold, though with different patterns. GraphCast outperformed ConvLSTM overall but had a sharper rise in exceedances (from 51.0% to 59.3%; $+\Delta 8.3\%$) compared to ConvLSTM's slower increase (from 72.0% to 75.2%; $+\Delta 3.2\%$). This suggests a notable limitation of GraphCast, especially in 2019 ($+\Delta 5.0\%$ interannual change): although it captures fine-scale features well, it is more sensitive to initial errors. Small inaccuracies, such as those from upwelling zones, can grow unpredictably over time, especially under strong mesoscale activity. This pattern

corresponds with observed upwelling trends in the California Current Upwelling System (CCUS) (Cropper et al., 2014), where summer upwelling has intensified in the permanent (21–26°N) and weak permanent (26–35°N) zones.

In terms of forecast quality, ConvLSTM reduces sensitivity to upwelling-driven variability by generating smoother outputs through anomaly averaging. This behavior is evident in the RA score: ConvLSTM tends to be underactive, while GraphCast produces noisier, overactive forecasts. However, this smoothing limits ConvLSTM’s ability to resolve fine-scale dynamics, making it less suitable for high-resolution applications. By contrast, the GNN captures spatial dependencies more effectively, achieving lower RMSE and higher ACC across most lead times. Yet, GraphCast’s higher standard deviation at later lead-times reveals sensitivity to local ocean variability, while ConvLSTM’s convolutional structure inherently suppresses small-scale fluctuations. This trade-off prioritizes stability over accuracy—ConvLSTM’s visually coherent forecasts lack the granularity to capture localized dynamics, whereas GraphCast’s graph-based connectivity preserves fine-scale patterns.

Integrating GraphCast fine-scale representation with ConvLSTM stabilizing smoothing could mitigate error propagation while retaining critical dynamic structures. Their divergent performances stem from fundamental architectural differences: GraphCast graph-based approach resolves ocean patterns more precisely, while ConvLSTM convolutions blur spatially nuanced features.

GraphCast demonstrates the transformative potential of MLOP, achieving 75.5% and 47% RMSE reductions compared to GLORYS at 5-day and 10-day lead times while generating 20-day SST forecasts in just 2.3 minutes. In contrast, GLORYS provides a comprehensive multi-variable ocean state representation, including salinity and velocity components, but requires 4 hours for a 7-day simulation. Despite its broad scope, GLORYS incorporates uncertainties stemming from data assimilation methods, parameterizations, and model biases. GraphCast’s exceptional speed makes it a powerful tool for real-time forecasting. Expanding its capabilities to include variables like salinity and currents is crucial for broader oceanographic applications.

On January 1, 2020, both MLOP models experienced a sharp drop in forecast skill, coinciding with a significant increase in satellite instrumental error. This highlights the strong dependence of data-driven ocean forecasts on the quality of initial conditions. As noted by Bouall  gue et al. (2024), initializing with higher-quality data, such as operational IFS analyses instead of ERA5, can significantly improve forecast accuracy. The substantial observational errors introduced at initialization likely degraded forecast performance, depending on each model’s sensitivity to error growth, where the ocean’s chaotic nature amplifies small discrepancies. This event underscores the vulnerability of machine learning models to sporadic observational anomalies and reinforces the need for robust quality-control procedures, improved error modeling, and advanced data assimilation strategies.

The seasonal analysis using a three-month moving window average shows that GraphCast consistently outperforms both ConvLSTM and the GLORYS reanalysis, although performance varies by season. The deep learning models exhibit their weakest performance during summer (JJA), when upwelling intensifies due to the seasonal shift of the Azores High and the Inter-Tropical Convergence Zone (ITCZ) (Wooster et al., 1976; Cropper et al., 2014). Despite these challenging conditions, GraphCast reduces RMSE by 28.4% compared to ConvLSTM and by 75.3% relative to GLORYS, suggesting that its graph-based architecture more effectively captures oceanographic processes linked to intensified upwelling.

In contrast, ConvLSTM performs notably worse in JJA due to its convolutional architecture, which smooths the forecasts and filters out high-frequency variability. While this reduces noise, it also suppresses the sharp gradients and nonlinear interactions characteristic of active upwelling regimes, limiting its ability to resolve fine-scale summer dynamics. In winter (DJF), when upwelling is weaker and ocean conditions are more stable, this smoothing becomes less detrimental, and GraphCast’s advantage becomes even more pronounced, excelling at capturing steady-state patterns. During spring (MAM) and autumn (SON), both models face similar challenges, and their performance converges.

The spatial distribution of errors reveals high RMSE values near Cape Ghir, Bojador, and Blanco—areas characterized by intense upwelling, complex bathymetry, and coastal current interactions. These conditions pose significant challenges for both numerical and deep learning models, which still struggle to fully resolve filament generation processes driven by wind forcing, coastline irregularities, and mesoscale instabilities. The absence of explicit atmospheric forcings (e.g., wind stress) and bathymetric inputs in the deep learning framework likely limits its ability to disentangle these mechanisms. Additionally, the resolution of satellite-derived L4 data used for training cannot capture submesoscale features ($< 10\text{ km}$), potentially introducing systematic biases. This aligns with findings from recent studies showing that global climate models underestimate upwelling intensity due to coarse spatial resolution (Bindoff et al., 2019; Docquier et al., 2019), highlighting the need for regional downscaling and high-resolution ocean-atmosphere coupled models (Roberts et al., 2018).

A drawback of GraphCast is the emergence of triangular artifacts at longer lead times, visible as mesh-like patterns superimposed on SST predictions. Visual analysis links these artifacts to the decoder’s mesh structure, where each

triangular element updates SST values within its influence area. Since neighboring triangles share only one edge and two nodes, inconsistencies arise—especially in dynamic regions like coastal upwelling zones—where adjacent elements may experience different physical processes. This highlights the difficulty of ensuring coherence across mesh elements under heterogeneous ocean conditions.

Mitigating these artifacts presents both computational and architectural challenges. Refining the mesh can improve spatial consistency by reducing each triangle’s influence area, as seen in higher-resolution graph-based models (Lam et al., 2023; Oskarsson et al., 2023; Holmberg et al., 2024), but at significant computational cost. Alternatively, increasing node connectivity in the decoder, drawing information from more than three nodes using attention mechanisms, could smooth transitions across triangles mitigating artifacts by distributing updates across a broader mesh neighborhood. A third strategy involves integrating convolutional layers into the decoder to blend localized element outputs, leveraging their proven spatial smoothing effect. This approach, supported by ConvLSTM’s artifact-free outputs and lower variability, has been effective in reducing grid artifacts in transformer-based models (Couairon et al., 2024). Future work should explore hybrid architectures that combine adaptive mesh connectivity with convolutional smoothing, preserving GraphCast’s ability to resolve large-scale dynamics while enhancing robustness to localized errors.

GraphCast outperforms ConvLSTM and GLORYS in coastal upwelling zones, reducing errors by up to 28.5% in upwelling regions, highlighting its superior ability to capture atmosphere–ocean interactions and fine-scale spatial gradients. However, persistent errors in these areas show that neither numerical nor deep learning models achieve optimal performance without targeted improvements. Despite these challenges, GraphCast’s strong SST forecasting skill demonstrates the potential of machine learning for ocean prediction. Advancing toward operational viability will require reducing sensitivity to initial condition errors—especially under strong mesoscale activity—through improved data assimilation or hybrid modeling strategies. Architectural enhancements such as convolutional layers or attention mechanisms may help mitigate artifacts efficiently, while scaling to multivariate forecasts will demand balancing complexity and speed, possibly via latent representations and physics-informed regularization. Addressing regional and seasonal variability will also require adaptive graph structures and real-time forcing assimilation, with focused validation in key regions like upwelling systems to ensure robustness and reliability.

6 Conclusion

This work presented a detailed evaluation of deep learning architectures for oceanographic forecasting, focusing on an adapted version of GraphCast applied to the Canary Upwelling System. GraphCast achieved high spatial resolution and forecasting accuracy, surpassing ConvLSTM and traditional reanalysis products like GLORYS. These results highlighted the capacity of graph-based models to resolve mesoscale features such as filaments and eddies below 20 km in scale, offering a promising direction for data-driven approaches in operational oceanography.

Nonetheless, our findings also revealed important trade-offs. While GraphCast excelled in spatial precision, it exhibited increased sensitivity to initial condition errors, especially during summer upwelling seasons when ocean variability is highest. In contrast, ConvLSTM provided greater temporal stability due to its convolutional smoothing, but at the cost of lower spatial fidelity. These differences underscore the potential of hybrid architectures that combine the strengths of both models to balance robustness and detail in multi-day forecasts.

The analysis further showed distinct seasonal and regional error patterns. Both models performed best under winter steady-state conditions and struggled during high-variability summer periods. Spatial error concentrations were particularly pronounced near coastal capes, where bathymetric complexity and submesoscale processes below the resolution of current L4 SST training data introduced persistent inaccuracies. Additionally, the dependence on SST limited the capacity of models to provide a full representation of the ocean state, highlighting the need for physics-informed, multivariate model extensions.

Future developments should focus on hybrid architectures, improved data assimilation methods, and curated high-resolution training datasets capable of capturing small-scale dynamics. Expanding the prediction scope beyond SST—incorporating currents, salinity, and other variables—will enhance model relevance for marine resource management. Regionally adaptive graph designs that account for bathymetric and dynamical physical gradients will also be essential for mitigating systematic errors in critical areas such as coastal capes.

This study supports the growing role of machine learning in ocean prediction, while pointing out key innovations to fully realize its potential in operational settings. Our findings position data-driven ocean prediction models as scalable alternatives to traditional numerical systems, with the ability to enhance the resolution, speed, and adaptability of next-generation ocean forecasting systems.

CRedit authorship contribution statement

Giovanny Alejandro C-Londoño: Writing – original draft, Writing – review & editing, Conceptualization, Formal analysis, Methodology, Investigation, Data curation, Software, Validation, Visualization. **Javier Sánchez:** Writing – original draft, Writing – review & editing, Project administration, Supervision, Conceptualization, Formal analysis, Methodology, Resources. **Ángel Rodríguez-Santana:** Writing – review & editing, Supervision, Conceptualization.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements

The authors would like to express their gratitude to Javier J. Sánchez-Medina from Centro de Innovación para la Sociedad de la Información (CICEI) at the Universidad de Las Palmas de Gran Canaria for providing access to their computing servers, which were essential for training the models. We also extend our sincere thanks to Mercator Ocean International (MOI) for sharing the PSY4V3R1 model forecasting data, which was provided upon request. Their support and collaboration were invaluable to the success of this research.

Declaration of Generative AI and AI-assisted technologies in the writing process

The authors used ChatGPT, Gemini, and Deepseek to improve the readability and language of the manuscript during the preparation of this work. After using these tools, the authors reviewed and edited the content as needed and took full responsibility for the content of the published article.

Data availability

The data and code used in this paper are available at the following URLs:

- The European North West Shelf/Iberia Biscay Irish Seas - High Resolution L4 Sea Surface Temperature Reprocessed is available at <https://doi.org/10.48670/moi-00153>
- GLORYS is available at <https://doi.org/10.48670/moi-00021>.
- The source code of our models is publicly available at <https://github.com/gacuervol/regional-graphcast-sst.git> under the MIT licence.

Computational library

The methodology is implemented in Python using the JAX library for efficient multi-GPU deep learning training. Atmospheric and oceanographic datasets are processed using Xarray, enabling scalable handling of multi-dimensional arrays, while evaluation metrics are computed with NumPy. The deep learning models are constructed using Google DeepMind Flax and Haiku frameworks, which are optimized for JAX, enabling high-performance model development.

A Model Configuration, training protocols and experimental setup

A.1 Training details

The dataset was split into training, validation, and test sets using an 80/10/10 ratio, based on years rather than individual days. This temporal division preserves the representativeness of seasonal cycles in each subset. The training set covers 1982 to 2012 (31 years), the validation set spans from 2013 to 2016 (4 years), and the test set includes 2017 to 2020 (4 years). This scheme prioritizes training with historical data while testing is conducted with the most recent information.

The samples in the training set are organized into windows of three consecutive dates, with each window overlapping the next one by a single day. This approach maintains continuity in prediction trajectories while maximizing the use of available data.

We used each set to generate time forcings, providing the model with temporal context for the current day and year. As in Lam et al. (2023), we calculated these forcings as seconds starting at the UNIX epoch from January 1, 1970, and encoded them using sine and cosine transformations to capture cyclical patterns.

For training, we organized the dataset into batches of eight samples. Each batch consists of three components: i) the input data, which includes the initial conditions for SST over two days; ii) the target data, representing the SST values for the lead times (days to predict); and iii) the forcing data, which includes the time-generated forcings for the lead times and the static land-sea mask. To ensure a stable optimization process, we applied shuffling exclusively to the training set, while the validation and test sets remained unshuffled to maintain their temporal structure.

We computed normalization factors, such as the mean, standard deviation, and standard deviation of temporal changes, using the training set. These factors were then applied to normalize the dataset, ensuring consistent scaling across all samples. This preprocessing step improved model stability and enhanced training performance by standardizing the input distribution.

We adopted the training methodologies from Lam et al. (2023), excluding the initial warm-up phase because our experiments did not indicate a significant impact. Our training strategy comprised a 150-epoch training phase, completed in approximately 64 hours. We adjusted the learning rate using a half-cosine decay function, gradually reducing it from 10^{-3} to zero and updating it after each iteration. The training phase was oriented to forecasting one lead time with samples of three-time instants, $(\mathbf{x}^{t-1}, \mathbf{x}^t) \rightarrow \mathbf{x}^{t+1}$.

We used gradient descent to optimize the loss function. For adaptive moment estimation and weight decay management, we utilized the AdamW optimizer with parameters $\beta_1 = 0.9$, $\beta_2 = 0.99$, $\epsilon = 10^{-8}$, and $\lambda = 0.1$. Furthermore, we implemented gradient clipping with a maximum norm value of 32 to ensure stable training dynamics.

A.2 Model configuration

The ConvLSTM model employed in this study was specifically designed for spatiotemporal ocean prediction. To ensure a fair comparison with GraphCast, the model was configured and trained under similar configurations. It comprised two stacked ConvLSTM layers: the first layer contained 8 feature channels, and the second layer reduced this to 1 feature channel. Both layers utilized a 3×3 kernel with a stride of 1 and padding of 1. The input sequences had a shape corresponding to a batch size of 8, two time steps, and spatial dimensions of 256×256 with a single-channel input. The model was optimized using the AdamW optimizer, with a learning rate of 10^{-2} and a weight decay of 0.1. Training was conducted over 150 epochs, with updates for single-time-step predictions.

The architecture of the graph-based approach employed three levels of multi-mesh refinement, M^3 , implemented by recursively subdividing triangular mesh elements into smaller components. The grid structure operated at a resolution of 0.05° derived from the satellite L4 product resolution in the training data. Spatial connectivity was regulated by a parameter that defined the neighborhood radius for each mesh node, set at 0.6 times the edge length. Within this range, grid nodes were connected to the corresponding mesh node in the encoder, ensuring localized information aggregation. The processor module performed 6 message-passing steps within an 8-dimensional MLP latent space, with all MLPs containing a single hidden layer. In the decoder, edge normalization for mesh-to-grid transitions was performed using the maximum edge length. This configuration effectively balanced model complexity and computational efficiency.

A.3 Spatially-weighted Loss function

Most weather prediction studies use the Mean Squared Error (MSE) to optimize the neural network parameters during training. Since the area of a spherical grid cell decreases toward the geographic poles, MSE is often adjusted by latitude to account for this distortion. Common approaches include weighting by the cosine of the latitude (Rasp et al., 2020; Lam et al., 2023) or by using the difference between the sines of the cell’s edges (Rasp et al., 2024), both of which serve as relative indicators of grid cell area. Without such adjustments, large errors in small high-latitude cells are treated the same as small errors in much larger equatorial cells, leading to spatial imbalance in the loss. Therefore, many global prediction studies use latitude-weighted MSE.

However, in this study, we adopt a different strategy that focuses the learning process exclusively on the ocean by applying a spatial mask derived from the structure of the data. This spatial mask assigns a weight of zero to grid cells over land and a weight of one to cells over the ocean. These binary weights are used to exclude land areas from contributing to the loss computation, effectively masking out irrelevant regions and ensuring that the loss is computed only over ocean areas. This approach is particularly well suited to our setup, which targets a geographically limited domain. In such cases, the variation in latitude across the domain is relatively small, and the differences in grid cell area are minimal. Therefore, the benefit of applying latitude-based weighting is negligible, and the masking strategy offers a more direct way to focus the training process on the regions of interest.

The loss function is thus defined as:

$$\mathcal{L}_{\text{MSE}} = \frac{1}{N_{\text{Rollout}}} \sum_{t=t_0+1}^{t_0+N_{\text{Rollout}}} \frac{1}{|\mathbb{G}_{\text{ocean}}|} \sum_{v_i \in \mathbb{G}_{\text{ocean}}} (\hat{x}_i^t - x_i^t)^2,$$

where

- t_0 is the initial forecast time,
- N_{Rollout} is the number of rollout steps used during training,
- v_i is a grid point defined by its latitude (ϕ) and longitude (λ),
- \mathbb{G} is the full spatial grid, and $|\mathbb{G}| = |\phi| \cdot |\lambda|$,
- $\mathbb{G}_{\text{ocean}} \subset \mathbb{G}$ is the subset of grid points located over the ocean, determined using a binary land–ocean mask $m_i \in \{0, 1\}$, where $m_i = 1$ indicates ocean and $m_i = 0$ indicates land.

A.4 Experimental setup

We used two computers to conduct our experiments. The first machine, a local workstation equipped with a single Quadro RTX 4000 GPU with 8 GB of memory, was used for hyperparameter tuning. The second machine, a remote server with 8 Quadro RTX 4000 GPUs and 8 GB of memory, was utilized for training the final model with the full dataset. The training phase on the remote server took approximately 64 hours. Our software stack included JAX, Haiku, Jraph, Optax, Jaxline, and xarray (Hoyer and Hamman, 2017) for model customization and training.

During the development phase on the local machine, we conducted a hyperparameter search to explore various configurations. Specifically, we investigated the number of message-passing steps $\{2, 6, 7\}$, latent size $\{2, 4, 8, 16\}$, and the mesh refinement level $\{2, 4, 6\}$ denoted as M^r . We tested batch sizes of 8 and 16, as preliminary results indicated that batch size had a negligible impact on the model’s performance. To efficiently cover a wide range of configurations, we performed a grid search over these three hyperparameters, limiting the number of combinations while ensuring a comprehensive exploration.

Once the hyperparameters were determined, we proceeded to train the final model on the remote server using the full dataset. To optimize the training process, we parallelized it using the Distributed Data-Parallel (DDP) strategy (Li et al., 2020). In this setup, each of the 8 GPUs held its copy of the model and processed a separate batch of data. After processing, we used the all-reduce operation to combine the gradients from all GPUs into a unified gradient, which was then used to update the model weights.

To accommodate the parallel processing, we added a new dimension to the data called *device*, with a size equal to 8 GPUs. This modified the data dimensions to 8 devices, 8 batches, 1 to n lead-time predictions, and 300×300 latitude-longitude size. Each GPU was assigned identical initialized weights, and during each training step, a forward pass was performed to calculate the loss and gradients for its assigned batch. At the end of the step, the gradients were summed on the CPU, and the resulting unified gradient was used to update the model weights across all GPUs.

This process was repeated iteratively until only residual batches remained. Since the number of residual batches was smaller than 8 (the number of GPUs), they could not be parallelized and were instead processed sequentially on the CPU.

B Graph Neural Network Model

Let $\mathbf{x} : \mathbb{R}^3 \rightarrow \mathbb{R}$ be a spatiotemporal variable containing the SST volume for our period of study. Each sequence value is represented as x_i^t , with i standing for node in a grid $v_i \in \mathcal{V}^g$, and t the time instant. For simplicity, let \mathbf{x}^t be the SST bidimensional array for time t and \mathbf{x}_i the temporal evolution of the SST in node i . Each node v_i has a latitude-longitude position given by (ϕ_i, λ_i) .

Our method is based on the GraphCast model (Lam et al., 2023), which is defined as an autoregressive model:

$$\hat{\mathbf{x}}^{t+1} = f(\mathbf{x}^t, \mathbf{x}^{t-1}), \quad (3)$$

where the $\hat{\mathbf{x}}^{t+1}$ map is a forecast obtained from two previous time instants. This output can be iteratively fed into the model to predict consecutive forecastings in a roll-out manner.

The graph of the model, $\mathcal{G}(\mathcal{V}^g, \mathcal{V}^m, \mathcal{E}^m, \mathcal{E}^{g2m}, \mathcal{E}^{m2g})$, is composed of grid nodes, \mathcal{V}^g , mesh nodes, \mathcal{V}^m , bidirectional edges connecting mesh nodes, \mathcal{E}^m , and directed edges from grid to mesh nodes, \mathcal{E}^{g2m} , and vice-versa, \mathcal{E}^{m2g} .

Each node in the grid, $\mathbf{v}_i^g \in \mathcal{V}^g$, is built from the SST values in times t and $t - 1$, several external forcings in times $t - 1$, t and $t + 1$, and constant values, i.e., $\mathbf{v}_i^g = [x_i^t, x_i^{t-1}, \mathbf{f}_i^{t-1}, \mathbf{f}_i^t, \mathbf{f}_i^{t+1}, \mathbf{c}_i]$, in a given latitude-longitude position. The forcings are defined by $\mathbf{f} = [\sin(h), \cos(h), \sin(y), \cos(y)]$, with h the local time of day and y the year progress, normalized to $[0, 1)$. The constants are defined as $\mathbf{c}_i = [m_i^{0/1}, \cos(\phi_i), \sin(\lambda_i), \cos(\lambda_i)]$, with $m^{0/1}$ a binary land-sea mask and no physical forcings.

Each node in the mesh, $\mathbf{v}_i^m \in \mathcal{V}^m$, is defined as $\mathbf{v}_i^m = [\cos(\phi_i), \sin(\lambda_i), \cos(\lambda_i)]$ and the mesh edges, $\mathbf{e}_{s,r}^m \in \mathcal{E}^m$, from the sender node s to the receiver node r , contains the following features: $\mathbf{e}_{s,r}^m = [\text{distance}(s, r), \mathbf{p}_s - \mathbf{p}_r]$, i.e. the edge length and the difference between the 3D spatial locations of nodes s and r . Figure 2 shows a representation of the multi-mesh with three levels of resolution.

Unidirectional edges from the grid to the mesh are similarly defined by $\mathbf{e}_{s,r}^{g2m} = [\text{distance}(s, r), \mathbf{p}_s - \mathbf{p}_r]$. In this case, the sender node is on the grid and the receiver is on the mesh. An edge is added between two nodes if $\text{distance}(s, r)$ is smaller than or equal to 0.6 times the length of the edges at the finest scale.

Similarly, unidirectional edges from the mesh to the grid are defined by $\mathbf{e}_{s,r}^{m2g} = [\text{distance}(s, r), \mathbf{p}_s - \mathbf{p}_r]$. An edge is created with the three mesh nodes of the triangular face that contains the grid node.

B.1 Encoder

The *encoder* maps the input data from the latitude-longitude grid into the multi-scale mesh. It relies on multi-layer perceptrons (MLP) to embed the graph variables, \mathbf{v}_i^g , \mathbf{v}_i^m , $\mathbf{e}_{s,r}^m$, $\mathbf{e}_{s,r}^{g2m}$ and $\mathbf{e}_{s,r}^{m2g}$, into a latent space:

$$\begin{aligned}\tilde{\mathbf{v}}_i^g &= \text{MLP}_{\mathcal{V}^g}(\mathbf{v}_i^g), \\ \tilde{\mathbf{v}}_i^m &= \text{MLP}_{\mathcal{V}^m}(\mathbf{v}_i^m), \\ \tilde{\mathbf{e}}_{s,r}^m &= \text{MLP}_{\mathcal{E}^m}(\mathbf{e}_{s,r}^m), \\ \tilde{\mathbf{e}}_{s,r}^{g2m} &= \text{MLP}_{\mathcal{E}^{g2m}}(\mathbf{e}_{s,r}^{g2m}), \\ \tilde{\mathbf{e}}_{s,r}^{m2g} &= \text{MLP}_{\mathcal{E}^{m2g}}(\mathbf{e}_{s,r}^{m2g}).\end{aligned}\tag{4}$$

Then, the information from the grid is transferred to the mesh using interaction networks (IN) (Battaglia et al., 2016). First, the edges are updated with information from the sending and receiving nodes,

$$\mathbf{d}\tilde{\mathbf{e}}_{s,r}^{g2m} = \text{MLP}_{\mathcal{E}^{g2m}}([\tilde{\mathbf{e}}_{s,r}^{g2m}, \mathbf{s}^g, \mathbf{r}^m]),\tag{5}$$

and the mesh nodes are updated by aggregating the edges as

$$\mathbf{d}\tilde{\mathbf{v}}_i^m = \text{MLP}_{\mathcal{V}^m}\left([\tilde{\mathbf{v}}_i^m, \sum_{s \in \mathcal{V}^g; r = \tilde{\mathbf{v}}_i^m} \tilde{\mathbf{e}}_{s,r}^{g2m}]\right),\tag{6}$$

with s the set of nodes in the grid that have an edge towards $\tilde{\mathbf{v}}_i^m$. The grid nodes are also updated as

$$\mathbf{d}\tilde{\mathbf{v}}_i^g = \text{MLP}_{\mathcal{V}^g}(\tilde{\mathbf{v}}_i^g),\tag{7}$$

All MLPs in these equations are independent and do not share their parameters. Finally, the latent variables are updated with residual connections as

$$\begin{aligned}\tilde{\mathbf{v}}_i^g &\leftarrow \tilde{\mathbf{v}}_i^g + \mathbf{d}\tilde{\mathbf{v}}_i^g, \\ \tilde{\mathbf{v}}_i^m &\leftarrow \tilde{\mathbf{v}}_i^m + \mathbf{d}\tilde{\mathbf{v}}_i^m, \\ \tilde{\mathbf{e}}_{s,r}^{g2m} &\leftarrow \tilde{\mathbf{e}}_{s,r}^{g2m} + \mathbf{d}\tilde{\mathbf{e}}_{s,r}^{g2m}.\end{aligned}\tag{8}$$

B.2 Processor

The *processor* performs learned message-passing through several layers on the multi-mesh. First, the edges are updated through an MLP, concatenating the mesh edge with the receiver and sender nodes' latent variables as

$$\mathbf{d}\tilde{\mathbf{e}}_{s,r}^m = \text{MLP}_{\mathcal{E}^m}([\tilde{\mathbf{e}}_{s,r}^m, \mathbf{s}^m, \mathbf{r}^m]).\tag{9}$$

Mesh nodes are updated by aggregating the edges as

$$\mathbf{d}\tilde{\mathbf{v}}_i^m = \text{MLP}_{\mathcal{V}^m}\left([\tilde{\mathbf{v}}_i^m, \sum_{s \in \mathcal{V}^m; r = \tilde{\mathbf{v}}_i^m} \tilde{\mathbf{e}}_{s,r}^m]\right),\tag{10}$$

and the variables are updated as a residual connection as

$$\begin{aligned}\tilde{\mathbf{v}}_i^m &\leftarrow \tilde{\mathbf{v}}_i^m + \mathbf{d}\tilde{\mathbf{v}}_i^m, \\ \tilde{\mathbf{e}}_{s,r}^m &\leftarrow \tilde{\mathbf{e}}_{s,r}^m + \mathbf{d}\tilde{\mathbf{e}}_{s,r}^m.\end{aligned}\quad (11)$$

This process represents one layer of the processor. It contains multiple iterative layers with independent MLP parameters that perform several message-passing operations.

B.3 Decoder

The *decoder* performs the inverse mapping from the mesh to the latitude-longitude grid. First, the edges from the mesh to the grid are updated as

$$\mathbf{d}\tilde{\mathbf{e}}_{s,r}^{m2g} = \text{MLP}_{\mathcal{E}^{m2g}}[\tilde{\mathbf{e}}_{s,r}^{m2g}, \mathbf{s}^m, \mathbf{r}^g]. \quad (12)$$

The grid nodes are then updated, aggregating the information of the three edges that arrive at the grid node:

$$\mathbf{d}\tilde{\mathbf{v}}_i^g = \text{MLP}_{\mathcal{V}^g} \left([\tilde{\mathbf{v}}_i^g, \sum_{s \in \mathcal{V}^g; r = \tilde{\mathbf{v}}_i^g} \mathbf{d}\tilde{\mathbf{e}}_{s,r}^{m2g}] \right). \quad (13)$$

A residual connection is used to update the information of the grid nodes coming from the embedding in the encoder:

$$\tilde{\mathbf{v}}_i^g \leftarrow \tilde{\mathbf{v}}_i^g + \mathbf{d}\tilde{\mathbf{v}}_i^g, \quad (14)$$

and the output prediction is obtained with another MLP as

$$\hat{\mathbf{y}}^t = \text{MLP}_{\mathcal{V}^g}(\tilde{\mathbf{v}}_i^g). \quad (15)$$

Finally, the forecasting is obtained as a residual connection with the input data as

$$\hat{\mathbf{x}}^{t+1} = \mathbf{x}^t + \hat{\mathbf{y}}^t$$

All input variables are normalized to zero mean and unit variance. The output variable $\hat{\mathbf{y}}^t$ is multiplied by the average standard deviation of the temporal change $\mathbf{d}\hat{\mathbf{x}} = \hat{\mathbf{x}}^{t+1} - \hat{\mathbf{x}}^t$ computed from the train set.

All MLPs have one hidden layer with a latent dimension of 8, the same as the output layer. The size of the output layer in the last MLP, corresponding to the *decoder*, is one for the prediction of the value of the SST at each node.

C Description of the evaluation metrics

In this study, we used four metrics to assess the performance of the ocean prediction models: the Anomaly Correlation Coefficient (ACC), the Root Mean Square Error (RMSE), the Activity (ACT), and the Bias. While there are several ways to define those metrics, we use the approach outlined by Bouallègue et al. (2024) as it aligns with the scores set by the World Meteorological Organization (WMO) and the ECMWF. Below is a detailed description of each metric, including its mathematical formulation and interpretation.

The RMSE quantifies the average magnitude of the prediction error, providing a measure of the overall accuracy of the model. The ACC evaluates the linear association between the predicted and observed anomalies, serving as an indicator of the model's skill in capturing the spatiotemporal variability of the ocean field. The Bias represents the systematic offset between predictions and observations, while ACT is defined as the standard deviation of the predicted anomalies and reflects the model's ability to reproduce the observed variability amplitude given a measure of the smoothness of the forecast.

C.1 Anomaly Correlation Coefficient (ACC)

The ACC measures the correlation between the predicted and observed anomalies, evaluating the model's ability to predict deviations from the climatology. It is calculated as:

$$\text{ACC} = \frac{(\overline{a_f - \bar{a}_f})(\overline{a_o - \bar{a}_o})}{\sqrt{(\overline{a_f - \bar{a}_f})^2} \sqrt{(\overline{a_o - \bar{a}_o})^2}}, \quad (16)$$

where $a_f = \hat{x}_i^t - c_i^t$ denotes the predicted anomaly (prediction deviation from climatology c_i^t) and $a_o = x_i^t - c_i^t$ the observed anomaly (ground truth deviation) for a given time and location. $\overline{(\cdot)} = \frac{1}{|\mathbb{G}|} \sum_{v_i \in \mathbb{G}} w_\phi(\cdot)$ represents the latitudinal weighted spatial average, where w_ϕ is the latitude-weighting factor based on the cosine of the latitude expressed in radians (Geer, 2016; Rasp et al., 2020; Lam et al., 2023). This weighting scheme ensures that regions near the equator, where the longitudinal distance between grid points is larger, are not underrepresented in the score calculation.

The ACC ranges between -1 and 1, where values close to 1 indicate a high positive correlation between the predicted and observed anomalies. This metric is widely used in forecast verification studies (ECMWF).

C.2 Root Mean Square Error (RMSE)

The RMSE quantifies the average magnitude of the difference between the predicted and observed values, providing a measure of the model’s overall accuracy. It is defined as:

$$\text{RMSE} = \sqrt{(\hat{x}_i^t - x_i^t)^2}, \quad (17)$$

where \hat{x}_i^t is the predicted value and x_i^t is the ground truth value at a given time and grid cell.

A lower RMSE indicates better predictive accuracy, making it one of the most common metrics for forecast verification.

C.3 Activity (ACT)

The Activity assesses the model’s ability to reproduce the observed variability by measuring the standard deviation of the predicted and observed anomaly fields:

$$\text{ACT}_f = \sqrt{(a_f - \bar{a}_f)^2}, \quad \text{ACT}_o = \sqrt{(a_o - \bar{a}_o)^2}. \quad (18)$$

$$(19)$$

A similar ACT between predictions and observations suggests that the model accurately captures the spatial variability of the ocean field. This score was originally proposed by Thorpe et al. (2013) and later adopted by Bouallègue et al. (2024) to assess the smoothness of the forecast.

The relative activity (RA) of a model is defined as the ratio between the forecast ACT and the observed ACT (Bechtold et al., 2008). Analyzing RA as a function of forecast lead time reveals whether the model underestimates variability (producing smoother forecasts) or overestimates it (producing noisier forecasts):

$$\text{RA} = \frac{\text{ACT}_f}{\text{ACT}_o}. \quad (20)$$

C.4 Bias

The Bias quantifies the systematic difference between the predicted and observed values, indicating whether the model tends to overestimate or underestimate the observations. It is calculated as:

$$\text{Bias} = (\hat{x}_i^t - x_i^t). \quad (21)$$

A Bias close to zero implies that the model does not present systematic overestimation or underestimation of the observations. This metric is commonly used in the verification of ocean and atmospheric models.

These metrics provide a comprehensive assessment of the models’ performance, allowing us to quantify their predictive skill across different aspects of the forecast quality.

References

- Adcroft, A., Hallberg, R., Dunne, J. P., Samuels, B. L., Galt, J. A., Barker, C. H., and Payton, D. (2010). Simulations of underwater plumes of dissolved oil in the Gulf of Mexico. *Geophysical Research Letters*, 37(18).
- Arístegui, J., Barton, E. D., Álvarez Salgado, X. A., Santos, A. M. P., Figueiras, F. G., Kifani, S., Hernández-León, S., Mason, E., Machú, E., and Demarcq, H. (2009). Sub-regional ecosystem variability in the Canary Current upwelling. *Progress in Oceanography*, 83(1):33–48. Eastern Boundary Upwelling Ecosystems: Integrative and Comparative Approaches.

- Arístegui, J., Sangrá, P., Hernández-León, S., Cantón, M., Hernández-Guerra, A., and Kerling, J. (1994). Island-induced eddies in the Canary Islands. *Deep Sea Research Part I: Oceanographic Research Papers*, 41(10):1509–1525.
- Barton, E., Arístegui, J., Tett, P., Cantón, M., García-Braun, J., Hernández-León, S., Nykjaer, L., Almeida, C., Almunia, J., Ballesteros, S., Basterretxea, G., Escáñez, J., García-Weill, L., Hernández-Guerra, A., López-Laatzén, F., Molina, R., Montero, M., Navarro-Pérez, E., Rodríguez, J., van Lenning, K., Vélez, H., and Wild, K. (1998). The transition zone of the Canary Current upwelling region. *Progress in Oceanography*, 41(4):455–504.
- Battaglia, P. W., Hamrick, J. B., Bapst, V., Sánchez-González, A., Zambaldi, V., Malinowski, M., Tacchetti, A., Raposo, D., Santoro, A., Faulkner, R., Gulcehre, C., Song, F., Ballard, A., Gilmer, J., Dahl, G., Vaswani, A., Allen, K., Nash, C., Langston, V., Dyer, C., Heess, N., Wierstra, D., Kohli, P., Botvinick, M., Vinyals, O., Li, Y., and Pascanu, R. (2018). Relational inductive biases, deep learning, and graph networks. *arXiv e-prints*, page arXiv:1806.01261.
- Battaglia, P. W., Pascanu, R., Lai, M., Rezende, D., and Kavukcuoglu, K. (2016). Interaction Networks for learning about objects, relations and physics. Technical report, Cornell University.
- Bechtold, P., Köhler, M., Jung, T., Leutbecher, M., Rodwell, M., and Vitart, F. (2008). Advances in simulating atmospheric variability with IFS cycle 32r3. *ECMWF Newsletter No. 114 - Winter 2007/08*, 114:29–38.
- Bell, M. J., Lefebvre, M., Traon, P.-Y. L., Smith, N., and Wilmer-Becker, K. (2009). GODAE: the global ocean data assimilation experiment. *Oceanography*, 22(3):14–21.
- Bengio, Y., Simard, P., and Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, 5(2):157–166.
- Bi, K., Xie, L., Zhang, H., Chen, X., Gu, X., and Tian, Q. (2023). Accurate medium-range global weather forecasting with 3d neural networks. *Nature*, 619(7970):533–538.
- Bindoff, N. L., Cheung, W. W. L., Kairo, J. G., Arístegui, J., Guinder, V. A., Hallberg, R., Hilmi, N., Jiao, N., Karim, M. S., Levin, L., O’Donoghue, S., Cuicapusa, S. R. P., Rinkevich, B., Suga, T., Tagliabue, A., and Williamson, P. (2019). Changing ocean, marine ecosystems, and dependent communities. In Pörtner, H.-O., Roberts, D. C., Masson-Delmotte, V., Zhai, P., Tignor, M., Poloczanska, E., Mintenbeck, K., Alegría, A., Nicolai, M., Okem, A., Petzold, J., Rama, B., and Weyer, N. M., editors, *IPCC Special Report on the Ocean and Cryosphere in a Changing Climate*, pages 447–587. Cambridge University Press, Cambridge, UK and New York, NY, USA.
- Bodnar, C., Bruinsma, W. P., Lucic, A., Stanley, M., Brandstetter, J., Garvan, P., Riechert, M., Weyn, J., Dong, H., Vaughan, A., Gupta, J. K., Tambiratnam, K., Archibald, A., Heider, E., Welling, M., Turner, R. E., and Perdikaris, P. (2024). A foundation model of the atmosphere. Technical report, Cornell University.
- Bouallègue, Z. B., Clare, M. C. A., Magnusson, L., Gascón, E., Maier-Gerber, M., Janoušek, M., Rodwell, M., Pinault, F., Dramsch, J. S., Lang, S. T. K., Raoult, B., Rabier, F., Chevallier, M., Sandu, I., Dueben, P., Chantry, M., and Pappenberger, F. (2024). The rise of data-driven weather forecasting: A first statistical assessment of machine learning-based weather forecasts in an operational-like context. *Bulletin of the American Meteorological Society*, 105(6):E864 – E883.
- Brasseur, P. and Verron, J. (2006). The SEEK filter method for data assimilation in oceanography: a synthesis. *Ocean Dynamics*, 56:650–661.
- Cabanes, C., Grouazel, A., von Schuckmann, K., Hamon, M., Turpin, V., Coatanoan, C., Paris, F., Guinehut, S., Boone, C., Ferry, N., de Boyer Montégut, C., Carval, T., Reverdin, G., Pouliquen, S., and Le Traon, P.-Y. (2013). The CORA dataset: validation and diagnostics of in-situ ocean temperature and salinity measurements. *Ocean Science*, 9(1):1–18.
- Chambers, D. P., Cazenave, A., Champollion, N., Dieng, H., Llovel, W., Forsberg, R., von Schuckmann, K., and Wada, Y. (2017). Evaluation of the global mean sea level budget between 1993 and 2014. *Integrative study of the mean sea level and its components*, pages 315–333.
- Chattopadhyay, A., Gray, M., Wu, T., Lowe, A. B., and He, R. (2024). OceanNet: a principled neural operator-based digital twin for regional oceans. *Scientific Reports*, 14(1):21181.
- Chattopadhyay, A., Mustafa, M., Hassanzadeh, P., and Kashinath, K. (2021). Deep spatial transformers for autoregressive data-driven forecasting of geophysical turbulence. In *Proceedings of the 10th International Conference on Climate Informatics*, CI2020, page 106–112, New York, NY, USA. Association for Computing Machinery.
- Chattopadhyay, A., Sun, Y. Q., and Hassanzadeh, P. (2023). Challenges of learning multi-scale dynamics with AI weather models: Implications for stability and one solution. *arXiv e-prints*, page arXiv:2304.07029.
- CMEMS (2024). European North West Shelf/Iberia Biscay Irish seas - high resolution 14 sea surface temperature reprocessed.

- Couairon, G., Lessig, C., Charantonis, A., and Monteleoni, C. (2024). ArchesWeather: An efficient AI weather forecasting model at 1.5° resolution. Technical Report 2405.14527, Cornell University.
- Cropper, T. E., Hanna, E., and Bigg, G. R. (2014). Spatial and temporal seasonal trends in coastal upwelling off Northwest Africa, 1981–2012. *Deep Sea Research Part I: Oceanographic Research Papers*, 86:94–111.
- Cushman-Roisin, B., Tang, B., and Chassignet, E. P. (1990). Westward motion of mesoscale eddies. *Journal of Physical Oceanography*, 20(5):758 – 768.
- Dee, D. P., Uppala, S. M., Simmons, A. J., Berrisford, P., Poli, P., Kobayashi, S., Andrae, U., Balmaseda, M. A., Balsamo, G., Bauer, P., Bechtold, P., Beljaars, A. C. M., van de Berg, L., Bidlot, J., Bormann, N., Delsol, C., Dragani, R., Fuentes, M., Geer, A. J., Haimberger, L., Healy, S. B., Hersbach, H., Hólm, E. V., Isaksen, I., Kållberg, P., Köhler, M., Matricardi, M., McNally, A. P., Monge-Sanz, B. M., Morcrette, J.-J., Park, B.-K., Peubey, C., de Rosnay, P., Tavolato, C., Thépaut, J.-N., and Vitart, F. (2011). The ERA-Interim reanalysis: configuration and performance of the data assimilation system. *Quarterly Journal of the Royal Meteorological Society*, 137(656):553–597.
- Docquier, D., Grist, J. P., Roberts, M. J., Roberts, C. D., Semmler, T., Ponsoni, L., Massonnet, F., Sidorenko, D., Sein, D. V., Iovino, D., et al. (2019). Impact of model resolution on Arctic sea ice and North Atlantic ocean heat transport. *Climate Dynamics*, 53:4989–5017.
- Dueben, P. D. and Bauer, P. (2018). Challenges and design choices for global weather and climate models based on machine learning. *Geoscientific Model Development*, 11(10):3999–4009.
- Estrada-Allis, S. N., Rodríguez-Santana, Á., Naveira-Garabato, A. C., García-Weil, L., Arcos-Pulido, M., and Emelianov, M. (2023). Enhancement of turbulence and nutrient fluxes within an Eastern Boundary Upwelling Filament: a diapycnal entrainment approach. *Frontiers in Marine Science*, Volume 10 - 2023.
- Ezraty, R., Girard-Ardhuin, F., Piollé, J.-F., Kaleschke, L., and Heygster, G. (2007). Arctic and Antarctic sea ice concentration and arctic sea ice drift estimated from special sensor microwave data. *Département d’Océanographie Physique et Spatiale, IFREMER, Brest, France and University of Bremen Germany*, 2.
- Falkowski, P. G., Ziemann, D., Kolber, Z., and Bienfang, P. K. (1991). Role of eddy pumping in enhancing primary production in the ocean. *Nature*, 352(6330):55–58.
- Galloudec, O. L., Chune, S. L., Nouel, L., Fernandez, E., Derval, C., Tressol, M., Dussurget, R., Biardeau, A., and Tonani, M. (2024). *Product User Manual for Global Ocean Physics Analysis and Forecasting Product*. Copernicus Marine Environment Monitoring Service (CMEMS). CMEMS-GLO-PUM-001-024, Issue 2.1.
- Geer, A. J. (2016). Significance of changes in medium-range forecast scores. *Tellus A: Dynamic Meteorology and Oceanography*, 68(1):30229.
- Gers, F. and Schmidhuber, J. (2000). Recurrent nets that time and count. In *Proceedings of the IEEE-INNS-ENNS International Joint Conference on Neural Networks. IJCNN 2000. Neural Computing: New Challenges and Perspectives for the New Millennium*, volume 3, pages 189–194 vol.3.
- Gers, F., Schmidhuber, J., and Cummins, F. (1999). Learning to forget: continual prediction with LSTM. In *1999 Ninth International Conference on Artificial Neural Networks ICANN 99. (Conf. Publ. No. 470)*, volume 2, pages 850–855 vol.2.
- Greff, K., Srivastava, R. K., Koutník, J., Steunebrink, B. R., and Schmidhuber, J. (2017). LSTM: A search space odyssey. *IEEE Transactions on Neural Networks and Learning Systems*, 28(10):2222–2232.
- Hagen, E., Zulicke, C., and Feistel, R. (1996). Near-surface structures in the Cape Ghir filament off Morocco. *Oceanologica Acta*, 19(6):577–598.
- Hausmann, U. and Czaja, A. (2012). The observed signature of mesoscale eddies in sea surface temperature and the associated heat transport. *Deep Sea Research Part I: Oceanographic Research Papers*, 70:60–72.
- Hochreiter, S. and Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation*, 9(8):1735–1780.
- Holmberg, D., Clementi, E., and Roos, T. (2024). Regional Ocean Forecasting with Hierarchical Graph Neural Networks. *arXiv e-prints*, page arXiv:2410.11807.
- Hoyer, S. and Hamman, J. (2017). xarray: N-D labeled arrays and datasets in Python. *Journal of Open Research Software*, 5(1).
- Jean-Michel, L., Eric, G., Romain, B.-B., Gilles, G., Angélique, M., Marie, D., Clément, B., Mathieu, H., Olivier, L. G., Charly, R., Tony, C., Charles-Emmanuel, T., Florent, G., Giovanni, R., Mounir, B., Yann, D., and Pierre-Yves, L. T. (2021). The Copernicus global 1/12° oceanic and sea ice GLORYS12 reanalysis. *Frontiers in Earth Science*, 9.
- Keisler, R. (2022). Forecasting Global Weather with Graph Neural Networks. *arXiv e-prints*, page arXiv:2202.07575.

- Kochkov, D., Yuval, J., Langmore, I., Norgaard, P. C., Smith, J. A., Mooers, G., Klöwer, M., Lottes, J., Rasp, S., Düben, P. D., Hatfield, S., Battaglia, P. W., Sánchez-González, A., Willson, M., Brenner, M. P., and Hoyer, S. (2023). Neural general circulation models for weather and climate. *Nature*, 632:1060 – 1066.
- Lam, R., Sánchez-González, A., Willson, M., Wirsberger, P., Fortunato, M., Alet, F., Ravuri, S., Ewalds, T., Eaton-Rosen, Z., Hu, W., Merosse, A., Hoyer, S., Holland, G., Vinyals, O., Stott, J., Pritzel, A., Mohamed, S., and Battaglia, P. (2023). Learning skillful medium-range global weather forecasting. *Science*, 382(6677):1416–1421.
- Lang, S., Alexe, M., Chantry, M., Dramsch, J., Pinault, F., Raoult, B., Clare, M. C. A., Lessig, C., Maier-Gerber, M., Magnusson, L., Ben Bouallègue, Z., Prieto Nemesio, A., Dueben, P. D., Brown, A., Pappenberger, F., and Rabier, F. (2024). AIFS – ECMWF’s data-driven forecasting system. *arXiv e-prints*, page arXiv:2406.01465.
- Lellouche, J.-M., Greiner, E., Le Galloudec, O., Garric, G., Regnier, C., Drevillon, M., Benkiran, M., Testut, C.-E., Bourdalle-Badie, R., Gasparin, F., Hernandez, O., Levier, B., Drillet, Y., Remy, E., and Le Traon, P.-Y. (2018). Recent updates to the Copernicus Marine Service global ocean monitoring and forecasting real-time 1/12° high-resolution system. *Ocean Science*, 14(5):1093–1126.
- Lellouche, J.-M., Le Galloudec, O., Drévillon, 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. (2013). Evaluation of global monitoring and forecasting systems at Mercator Océan. *Ocean Science*, 9(1):57–81.
- Li, S., Zhao, Y., Varma, R., Salpekar, O., Noordhuis, P., Li, T., Paszke, A., Smith, J., Vaughan, B., Damania, P., and Chintala, S. (2020). PyTorch Distributed: Experiences on Accelerating Data Parallel Training. *arXiv e-prints*, page arXiv:2006.15704.
- Liu, Y., MacFadyen, A., Ji, Z.-G., and Weisberg, R. H. (2013). *Monitoring and modeling the deepwater horizon oil spill: a record breaking enterprise*. John Wiley & Sons.
- Madec, G., Bell, M., Benshila, R., Blaker, A., Boudrallé-Badie, R., Bricaud, C., Bruciaferri, D., Carneiro, D., Castrillo, M., Calvert, D., Chanut, J., Clementi, E., Coward, A., de Lavergne, C., Dobricic, S., Epicoco, I., Éthé, C., Fiedler, E., Ford, D., Furner, R., Ganderton, J., Graham, T., Harle, J., Hutchinson, K., Iovino, D., King, R., Lea, D., Levy, C., Lovato, T., Maisonnave, E., Mak, J., Sánchez, J. M. C., Martin, M., Martin, N., Martins, D., Masson, S., Mathiot, P., Mele, F., Mocavero, S., Moulin, A., Müller, S., Nurser, G., Oddo, P., Paronuzzi, S., Paul, J., Peltier, M., Person, R., Rousset, C., Rynders, S., Samson, G., Schroeder, D., Storkey, D., Storto, A., Téchené, S., Vancoppenolle, M., and Wilson, C. (2024). Nemo ocean engine reference manual.
- Madec, G. et al. (2008). NEMO ocean engine. note du Pôle de modélisation. *Institut Pierre-Simon Laplace (IPSL), France*, 27:1288–1619.
- Mourre, B., Aguiar, E., Juzà, M., Hernández-Lasheras, J., Reyes, E., Heslop, E., Escudier, R., Cutolo, E., Ruiz, S., Mason, E., Pascual, A., and Tintoré, J. (2018). Assessment of high-resolution regional ocean prediction systems using multi-platform observations: Illustrations in the Western Mediterranean Sea. *New Frontiers in Operational Oceanography*.
- Mu, B., Peng, C., Yuan, S., and Chen, L. (2019). ENSO forecasting over multiple time horizons using ConvLSTM network and rolling mechanism. In *2019 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8.
- Mu, B., Qin, B., and Yuan, S. (2021). ENSO-ASC 1.0.0: ENSO deep learning forecast model with a multivariate air–sea coupler. *Geoscientific Model Development*, 14(11):6977–6999.
- Oskarsson, J., Landelius, T., and Lindsten, F. (2023). Graph-based Neural Weather Prediction for Limited Area Modeling. *arXiv e-prints*, page arXiv:2309.17370.
- Pelegrí, J., Arístegui, J., Cana, L., González-Dávila, M., Hernández-Guerra, A., Hernández-León, S., Marrero-Díaz, A., Montero, M., Sangrà, P., and Santana-Casiano, M. (2005). Coupling between the open ocean and the coastal upwelling region off Northwest Africa: water recirculation and offshore pumping of organic matter. *Journal of Marine Systems*, 54(1):3–37. A general study of the Spanish North Atlantic boundaries: an interdisciplinary approach.
- Pfaff, T., Fortunato, M., Sánchez-González, A., and Battaglia, P. W. (2020). Learning Mesh-Based Simulation with Graph Networks. *arXiv e-prints*, page arXiv:2010.03409.
- Piollé, J.-F. and Autret, E. (2023). QUID for SST TAC Products SST_GLO_SST_L3S_NRT_OBSERVATIONS_010_010. Quality Information Document 2.1, E.U. Copernicus Marine Service Information (CMEMS).
- Price, I., Sánchez-González, A., Alet, F., Andersson, T. R., El-Kadi, A., Masters, D., Ewalds, T., Stott, J., Mohamed, S., Battaglia, P. W., Lam, R., and Willson, M. (2024). Probabilistic weather forecasting with machine learning. *Nature*, 637:84 – 90.
- Pujol, M.-I., Faugère, Y., Taburet, G., Dupuy, S., Pelloquin, C., Ablain, M., and Picot, N. (2016). DUACS DT2014: the new multi-mission altimeter data set reprocessed over 20 years. *Ocean Science*, 12(5):1067–1090.

- Rasp, S., Dueben, P. D., Scher, S., Weyn, J. A., Mouatadid, S., and Thuerey, N. (2020). WeatherBench: a benchmark data set for data-driven weather forecasting. *Journal of Advances in Modeling Earth Systems*, 12(11):e2020MS002203.
- Rasp, S., Hoyer, S., Merose, A., Langmore, I., Battaglia, P., Russel, T., Sánchez-González, A., Yang, V., Carver, R., Agrawal, S., Chantry, M., Bouallegue, Z. B., Dueben, P., Bromberg, C., Sisk, J., Barrington, L., Bell, A., and Sha, F. (2024). WeatherBench 2: A benchmark for the next generation of data-driven global weather models. Technical Report 2308.15560, Cornell University.
- Rio, M. H., Guinehut, S., and Larnicol, G. (2011). New CNES-CLS09 global mean dynamic topography computed from the combination of GRACE data, altimetry, and in situ measurements. *Journal of Geophysical Research: Oceans*, 116(C7).
- Roberts, M., Vidale, P., Senior, C., Hewitt, H., Bates, C., Berthou, S., Chang, P., Christensen, H., Danilov, S., Demory, M.-E., et al. (2018). The benefits of global high resolution for climate simulation: process understanding and the enabling of stakeholder decisions at the regional scale. *Bulletin of the American Meteorological Society*, 99(11):2341–2359.
- Sangrà, P., Pascual, A., Ángel Rodríguez-Santana, Machín, F., Mason, E., McWilliams, J. C., Pelegrí, J. L., Dong, C., Rubio, A., Arístegui, J., Ángeles Marrero-Díaz, Hernández-Guerra, A., Martínez-Marrero, A., and Auladell, M. (2009). The Canary Eddy Corridor: A major pathway for long-lived eddies in the subtropical North Atlantic. *Deep Sea Research Part I: Oceanographic Research Papers*, 56(12):2100–2114.
- Sangrà, P., Troupin, C., Barreiro-González, B., Desmond Barton, E., Orbi, A., and Arístegui, J. (2015). The cape Ghir filament system in August 2009 (NW Africa). *Journal of Geophysical Research: Oceans*, 120(6):4516–4533.
- Scher, S. (2018). Toward data-driven weather and climate forecasting: Approximating a simple general circulation model with deep learning. *Geophysical Research Letters*, 45(22):12,616–12,622.
- Shi, X., Chen, Z., Wang, H., Yeung, D.-Y., Wong, W.-k., and Woo, W.-c. (2015). Convolutional LSTM Network: A Machine Learning Approach for Precipitation Nowcasting. *arXiv e-prints*, page arXiv:1506.04214.
- Sommer, J. L., Chassignet, E. P., and Wallcraft, A. J. (2018). Ocean circulation modeling for operational oceanography: Current status and future challenges. *New Frontiers in Operational Oceanography*.
- Szekely, T., Gourrion, J., Pouliquen, S., Reverdin, G., and Merceur, F. (2019). CORA, coriolis ocean dataset for reanalysis. *Sea Scientific Open Data Publication*.
- Thompson, B., Sánchez, C., Heng, B. C. P., Kumar, R., Liu, J., Huang, X.-Y., and Tkalich, P. (2021). Development of a MetUM (v 11.1) and NEMO (v 3.6) coupled operational forecast model for the maritime continent – Part 1: Evaluation of ocean forecasts. *Geoscientific Model Development*, 14(2):1081–1100.
- Thorpe, A., Bauer, P., Magnusson, L., and Richardson, D. (2013). An evaluation of recent performance of ECMWF's forecasts. *ECMWF Newsletter No. 137 - autumn 2013*, pages 15–18.
- Treguier, A.-M., Chassignet, E. p., Le Boyer, A., and Pinardi, N. (2017). Modeling and forecasting the "weather of the ocean" at the mesoscale. *Journal Of Marine Research*, 75(3):301–329.
- Troupin, C., Mason, E., Beckers, J.-M., and Sangrà, P. (2012). Generation of the Cape Ghir upwelling filament: A numerical study. *Ocean Modelling*, 41:1–15.
- Wang, X., Wang, R., Hu, N., Wang, P., Huo, P., Wang, G., Wang, H., Wang, S., Zhu, J., Xu, J., Yin, J., Bao, S., Luo, C., Zu, Z., Han, Y., Zhang, W., Ren, K., Deng, K., and Song, J. (2024). XiHe: A Data-Driven Model for Global Ocean Eddy-Resolving Forecasting. Technical Report arXiv:2402.02995, Cornell University.
- Watters, N., Tacchetti, A., Weber, T., Pascanu, R., Battaglia, P., and Zoran, D. (2017). Visual Interaction Networks. *arXiv e-prints*, page arXiv:1706.01433.
- Wooster, W. S., Bakun, A., and McLain, D. R. (1976). The seasonal upwelling cycle along the eastern boundary of the North Atlantic. *Journal of Marine Research*, 34(2):131–141.
- Yang, J., Zhang, T., Zhang, J., Lin, X., Wang, H., and Feng, T. (2024). A ConvLSTM nearshore water level prediction model with integrated attention mechanism. *Frontiers in Marine Science*, 11.
- Yang, X., Zhang, F., Sun, P., Li, X., Du, Z., and Liu, R. (2022). A spatio-temporal graph-guided convolutional LSTM for tropical cyclones precipitation nowcasting. *Applied Soft Computing*, 124:109003.
- Zeng, X., Li, Y., and He, R. (2015). Predictability of the loop current variation and eddy shedding process in the Gulf of Mexico using an artificial neural network approach. *Journal of Atmospheric and Oceanic Technology*, 32(5):1098 – 1111.