# Uncertainty-permitting machine learning reveals sources of dynamic sea level predictability across daily-to-seasonal timescales

ANDREW BRETTIN[a] , LAURE ZANNA[a] , ELIZABETH A. BARNES[b]

[a] *Courant Institute of Mathematical Sciences, New York University, New York, New York,* [b] *Department of Atmospheric Science, Colorado State University, Fort Collins, Colorado*

ABSTRACT: Reliable dynamic sea level forecasts are hindered by numerous sources of uncertainty on daily-to-seasonal timescales (1–180 days) due to atmospheric boundary conditions and internal ocean variability. Studies have demonstrated that certain initial states can extend predictability horizons; thus, identifying these initial conditions may help improve forecast skill. Here, we identify sources of dynamic sea level predictability on daily-to-seasonal timescales using neural networks trained on CESM2 large ensemble data to forecast dynamic sea level. The forecasts yield not only a point estimate for sea level but also a standard deviation to quantify forecast uncertainty based on the initial conditions. Forecasted uncertainties can be leveraged to identify state-dependent sources of predictability at most locations and forecast leads. Network forecasts, particularly in the low-latitude Indo-Pacific, exhibit skillful deterministic predictions and skillfully forecast exceedance probabilities relative to local linear baselines. For networks trained at Guam and in the western Indian Ocean, the transfer of sources of predictability from local sources to remote sources is presented by the deteriorating utility of initial condition information for predicting exceedance events. Propagating Rossby waves are identified as a potential source of predictability for dynamic sea level at Guam. In the Indian Ocean, persistence of thermosteric sea level anomalies from the Indian Ocean Dipole may be a source of predictability on subseasonal timescales, but El Niño drives predictability on seasonal timescales. This work shows how uncertainty-quantifying machine learning can help identify changes in sources of state-dependent predictability over a range of forecast leads.

SIGNIFICANCE STATEMENT: Uncertainty-quantifying neural networks trained to forecast dynamic sea level anomalies on daily-to-seasonal timescales skillfully identify sources of state-dependent predictability and forecast exceedance probabilities. We use these predicted probabilities to identify how sources of state-dependent predictability change over daily-to-seasonal timescales. At Guam, Rossby waves are a source of dynamic sea level predictability. In the western Indian Ocean, El Niño emerges as a source of predictability, possibly by acting as a precursor to the Indian Ocean Dipole. This work may clarify how sources of predictability change between daily and seasonal timescales, which may help improve forecasts at these lead times.

## 1. Introduction

Dynamic sea level, defined as the height of the sea surface above the geoid (excluding the inverse-barometric imprint from atmospheric loading), is modulated by a variety of processes from the atmosphere and ocean (Gregory et al. 2019; Griffies et al. 2016). Dynamic sea level responds to atmospheric forcing through local processes, such as surface buoyancy fluxes (Hochet et al. 2024; Cabanes et al. 2006; Gill and Niller 1973) and Ekman pumping (Qu et al. 2022; Piecuch and Ponte 2011) and through remote processes such as time-dependent Sverdrup dynamics (Chen et al. 2023; Qiu 2002) and baroclinic response to wind

stress (Cabanes et al. 2006). Intrinsic ocean variability impacts dynamic sea level by mass redistribution (Fukumori et al. 1998), advection and diffusion of steric anomalies (Hochet et al. 2024) and barotropic and baroclinic instabilities (Penduff et al. 2010).

The wide variety of processes impacting dynamic sea level means that sea level fluctuates on a continuum of timescales. These variations can be significant: for instance, variations in dynamic sea level associated with the El Niño-Southern Oscillation can be as large as 20-30cm (Becker et al. 2012), comparable to the observed increase in global mean sea level due to contemporary anthropogenic climate change over the past century (Frederikse et al. 2020). Dynamic sea level variability on seasonal-to-interannual timescales (1-24 months, Jacox et al. 2020) has received considerable attention, and numerous studies have used observations and simulations to investigate the predictability of dynamic sea level on these time horizons (Balmaseda et al. 2024; Wang et al. 2023; Doi et al. 2020; Fraser et al. 2019; Miles et al. 2014). The interest in seasonal-to-interannual dynamic sea level variability may be driven, in part, by significant relationships with indices of climate variability on these timescales, such as the El Niño-Southern Oscillation, the Indian Ocean Dipole, the Southern Annular Mode and the North Atlantic Oscillation (Roberts et al. 2016; Chowdhury et al. 2007a; Miles et al. 2014; Aparna et al. 2012; Kenigson et al. 2018).

Studies have also highlighted the potential utility of sea level forecasts on subseasonal-to-seasonal timescales (15-

*Corresponding author*: Andrew Brettin, brettin@cims.nyu.edu

60 days, DeMott et al. 2021; Amaya et al. 2022; Arcodia et al. 2024). For instance, improved sea level forecasts on these timescales can help municipalities mitigate damage from nuisance "fair-weather" floods, which have occurred at increasing frequencies due to global mean sea level rise (Hino et al. 2019; Li et al. 2022). However, forecasts on subseasonal-to-seasonal timescales have generally been regarded as a challenging timescale for prediction in the earth system (Vitart et al. 2017; Mariotti et al. 2018; NASEM 2016). While the slow internal variability of the ocean compared to the atmosphere can provide some predictability, memory from the initial conditions of the atmosphere is typically lost beyond timescales of two weeks (Krishnamurthy 2019; Lorenz 1969).

Recent approaches that have allowed for the development of useful forecasts of geophysical conditions on subseasonal-to-seasonal timescales have focused on identifying specific initial conditions that can result in more skillful forecasts (Mariotti et al. 2020). Studies have established that certain initial conditions from the atmosphere or ocean can provide more predictability for the dynamics of geophysical fields than others (Christensen et al. 2020; Frame et al. 2013; Kalnay and Dalcher 1987). Thus, identifying such state-dependent sources of predictability can enable forecasts to be made on lead times that would normally not be considered (Albers and Newman 2019). Understanding sources of state-dependent predictability—and how these sources vary by forecast lead—can help bridge the gap between daily and seasonal forecasts.

Machine learning approaches, such as artificial neural networks (ANN), can help identify sources of state-dependent predictability directly from data. Mayer and Barnes (2021) showed how a classification artificial neural network trained to predict the sign of geopotential height anomalies could be used for a priori identification of skillful forecasts. In particular, activation functions can be leveraged to output class probabilities targeting predictable outcomes and their associated initial conditions. As an alternative method for estimating forecast uncertainty, Gordon and Barnes (2022) explored state-dependent predictability of sea surface temperatures on decadal timescales using a regression neural network trained on a Gaussian maximum-likelihood based loss function. The regression networks yield aleatoric uncertainty estimates for the prediction which can be directly applied to identify modes of variability associated with climate predictability, such as the Atlantic Meridional Variability and Interdecadal Pacific Oscillation.

Here, we apply a similar approach to Gordon and Barnes (2022) to investigate state-dependent sources of predictability of dynamic sea level anomalies on daily-to-seasonal timescales (1–180 days). We train regression neural networks to make probabilistic forecasts of sea level using simulated fields from the Community Earth System Model, version 2 Large Ensemble project (Danaba-

soglu et al. 2020; Rodgers et al. 2021). We examine how sources of state-dependent predictability change over daily-to-seasonal timescales, and identify these sources over different forecast leads.

## 2. Methods

### a. Data

We use simulated fields from the Community Earth System Model, version 2 (CESM2) Large Ensemble (LENS2) dataset (Danabasoglu et al. 2020; Rodgers et al. 2021). CESM2 is a fully-coupled earth system model run using a 1° nominal horizontal resolution in the ocean and atmosphere. The atmosphere is simulated using the finite-volume dynamical core of the Community Atmosphere Model, version 6 (CAM6, Lin and Rood 1997) and the ocean is simulated by solving the primitive equations employing the hydrostatic approximation using the Parallel Ocean Program, version 2 (POP2, Smith et al. 2010). The data is from the 250-year simulation period 1850–2100, with radiative forcing prescribed by the historical record from 1850–2015 and by the CMIP6 SSP370 forcing scenario from 2016–2100 (O'Neill et al. 2016). Ensemble simulations are initialized from a combination of different "macro-perturbations" of a preindustrial simulation state (based on different AMOC phases) and "micro-perturbations" (by adding minuscule noise to surface air temperature fields), as detailed in Rodgers et al. (2021). In addition to the different initialization procedures for different ensemble members in the LENS2 project, a further distinction between sets of ensemble members is made by imposing two different diagnostic surface forcing fields from biomass burning: in 50 of the 100 ensemble members, the CMIP6 protocol is followed; in the remaining 50 ensemble members, an 11-year running mean filter is applied to smooth large interannual variability in the forcing fields. We select nine of the smoothed biomass-burning ensemble members for this analysis, using three micro-perturbations for each of the three macro-perturbations that were available. Using different ensemble members allows us to identify climate drivers of sea level predictability which are robust under different climate forcings and internal variability. Because the ocean is simulated on a displaced-dipole grid, all oceanic variables are regridded to the uniform spherical-coordinate atmospheric grid using bilinear interpolation prior to any analysis.

We generate uncertainty-permitting forecasts for 5-day averaged dynamic sea level anomalies on the 1° grid at various daily-to-seasonal forecasting leads. For CMIP6 model evaluations (Eyring et al. 2016; Fox-Kemper et al. 2021), dynamic sea level is defined by the deviation of the sea surface height from the global mean, excluding the inverse-barometer contribution to sea surface height from atmospheric pressure loading (Griffies et al. 2014; Gregory et al. 2019; Wunsch and Stammer 1997). In CESM2, it

is computed using the implicit free-surface formulation of the barotropic equations from Dukowicz and Smith (1994) (Fasullo et al. 2020; Smith et al. 2010).

As model inputs, we use 5-day averaged dynamic sea level (ZOS, Griffies et al. 2016), sea surface temperatures (SST), and surface zonal and meridional wind fields (UAS and VAS, respectively) from 60°S to 60°N. These input fields are coarsened to 5° resolution to reduce the input dimensionality and help with understanding large-scale drivers, yielding a feature vector of 6,014 inputs.

To obtain anomalies, we detrend and deseasonalize all variables. Variables are detrended using a locally-fitted fifth-order polynomial computed over the full 250-year period for each ensemble member. The seasonal cycle is similarly removed by subtracting climatological daily averages at each grid point. Of the nine ensemble members used, seven are used for training (128,233 samples), and one is used for validation and testing, respectively (18,319 samples each). Values are standardized in time as a preprocessing step prior to training (LeCun et al. 2002). That is, for each location and field variable, the mean and standard deviation are taken over all samples in the training dataset and used to standardize the training, validation, and testing set.

### b. Machine learning framework

For each prediction location and forecast lead, we train two fully-connected regression artificial neural networks (ANN) to forecast dynamic sea level: one network predicts a point estimate for the forecast, and the second quantifies the uncertainty associated with the point estimate. Figure 1 illustrates the set-up of each network. Each network consumes the 6,014-dimensional vector of SST, ZOS, UAS, and VAS at every 5° gridpoint as input, and issues predictions at a single location on the 1° grid. To obtain spatial coverage, we train networks at every other gridpoint latitude and longitude between 60°S and 60°N (6,590 locations) for time lags of 10, 20, 60, and 120 days. Each network is trained using Pytorch (Paszke et al. 2019) using 1 CPU with 16 GB memory.

Given an input $x_i \in \mathbb{R}^{6,014}$ and model parameters $\boldsymbol{\theta}$, the first network outputs a prediction $\hat{\mu}_i = \hat{\mu}(x_i|\boldsymbol{\theta})$ for the target variable $y_i \in \mathbb{R}$. During training, parameters of the model are adjusted to optimize the Mean Square Error (MSE) loss function of the predicted values $\hat{\boldsymbol{\mu}} = (\hat{\mu}_1, \ldots, \hat{\mu}_N)$ given the targets $\boldsymbol{y} = (y_1, \ldots, y_N)$:

$$\mathcal{L}_{\text{MSE}}(\boldsymbol{y}, \hat{\boldsymbol{\mu}}) = \frac{1}{N} \sum_{i=1}^{N} (y_i - \hat{\mu}_i)^2. \quad (1)$$

The point estimate networks contain two hidden layers of 10 nodes each. Hidden layer nodes are equipped with rectified linear unit (ReLU) activation functions. Because the

feature vector is relatively large, we apply dropout regularization (Srivastava et al. 2014) with a dropout probability of 0.1 between the input layer and the first hidden layer. Parameters are updated using the Adam optimizer (Kingma and Ba 2014) with a learning rate of $10^{-5}$ using batches of 32 samples. Models are trained up to 100 epochs, but early stopping is applied so that training is halted if the validation MSE has not decreased for 10 epochs. Early stopping is used not only as a regularization techinque (Nakkiran et al. 2019) but also to reduce the computational cost of training several thousand neural networks.

After the model for predicting the point estimates for forecasted sea level is trained, we train a separate network using the residuals $\boldsymbol{y} - \hat{\boldsymbol{\mu}}$ as targets to quantify the uncertainty associated with the prediction (Nix and Weigend 1994). Given an input $x_i$ and model parameters $\boldsymbol{\theta}$, the uncertainty network outputs a predicted uncertainty as a log-variance, $\log \hat{\sigma}_i^2 = \log\left(\hat{\sigma}(x_i|\boldsymbol{\theta})^2\right)$. (The log-variance is used as an output in lieu of the variance to enforce the non-negativity of the predicted standard deviation, which can be easily computed from the log-variance via a bijective mapping.) In this network, the output uncertainty is optimized using the Gaussian negative log-likelihood loss function instead of the MSE. The Gaussian likelihood of a set of independent, identically distributed observations $\boldsymbol{y} = (y_1, \ldots, y_N)$ given parameters $\hat{\boldsymbol{\mu}} = (\hat{\mu}_1, \ldots, \hat{\mu}_N)$ and $\hat{\boldsymbol{\sigma}}^2 = (\hat{\sigma}_1^2, \ldots \hat{\sigma}_N^2)$ is given by

$$p(\boldsymbol{y}|\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\sigma}}^2) = \prod_{i=1}^{N} \frac{1}{\hat{\sigma}_i \sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{y_i - \hat{\mu}_i}{\hat{\sigma}_i}\right)^2\right] \quad (2)$$

The negative log-likelihood over all samples is given by

$$\begin{aligned} \mathcal{L}_{\text{NLL}}(\boldsymbol{y}|\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\sigma}}^2) &= -\log\left(p(\boldsymbol{y}|\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\sigma}}^2)\right) \\ &= \frac{1}{2} \sum_{i=1}^{N} \left[\log(\hat{\sigma}_i^2) + \left(\frac{y_i - \hat{\mu}_i}{\hat{\sigma}_i}\right)^2\right] + C \end{aligned} \quad (3)$$

where $C = \frac{N}{2}\log(2\pi)$ is an immaterial constant under optimization. During training, the model parameters are adjusted to optimize the predicted standard deviations $\hat{\boldsymbol{\sigma}}$ given the residuals $\boldsymbol{y} - \hat{\boldsymbol{\mu}}$.

The residuals $\boldsymbol{y} - \hat{\boldsymbol{\mu}}$ are further standardized prior to training the uncertainty network. The network architecture, optimizers, regularization, and training procedure are identical to that of the mean network, except that a learning rate of $10^{-6}$ is used for the uncertainty network to ensure convergence.

Predicting a mean and standard deviation means that predictions for dynamic sea level are formulated as Gaussian probability density functions, as shown in Figure 2a. Although the point estimate may deviate from the target value of sea level, the uncertainty-quantifying network aims to quantify the error in the forecast by a predicted standard
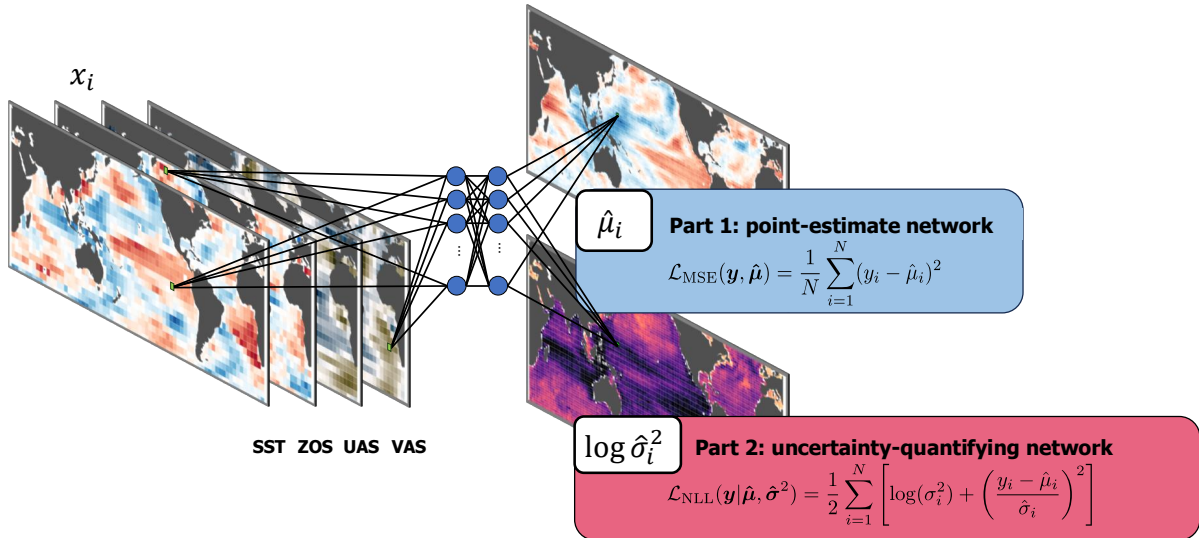
$x_i$

**Part 1: point-estimate network**
$\hat{\mu}_i$
$$\mathcal{L}_{\mathrm{MSE}}(\boldsymbol{y}, \hat{\boldsymbol{\mu}}) = \frac{1}{N}\sum_{i=1}^{N}(y_i - \hat{\mu}_i)^2$$

**SST ZOS UAS VAS**

$\log \hat{\sigma}_i^2$
**Part 2: uncertainty-quantifying network**
$$\mathcal{L}_{\mathrm{NLL}}(\boldsymbol{y}|\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\sigma}}^2) = \frac{1}{2}\sum_{i=1}^{N}\left[\log(\sigma_i^2) + \left(\frac{y_i - \hat{\mu}_i}{\hat{\sigma}_i}\right)^2\right]$$

FIG. 1: Network schematic. First, the point-estimate network consumes the input vector $x_i$ of 5°-coarsened sea surface temperatures, dynamic sea levels, and zonal and meridional winds and forecasts a point-estimate for dynamic sea level for a specified forecast lead at a given location on the 1° grid. This network is trained on the Mean Squared Error. Then, the uncertainty-quantifying network uses the same input fields to quantify the forecast uncertainty for a specified forecast lead at a given location on the 1° grid. This network is trained using the Gaussian negative log-likelihood loss. To obtain spatial coverage, separate networks are trained at 6,590 different gridpoints on the 1° grid for each forecast lead.

deviation $\hat{\sigma}$. This is shown in the scatterplot in Figure 2b, where predicted dynamic sea levels $\hat{\mu}$ and predicted uncertainties $\hat{\sigma}$ are plotted against the target dynamic sea level anomaly. While the point estimate typically gives an imperfect prediction for the dynamic sea level, the uncertainty-quantifying network leverages information from the initial conditions to help quantify the forecast uncertainty.

In general, it is not expected that the predicted standard deviation can be used to predict the residual $y_i - \hat{\mu}_i$, as the uncertainty-quantifying network uses the same input features as the point-estimate network. However, for well-calibrated forecasts, the predicted standard deviation will quantify the residual *in a statistical sense* when averaged over many samples. For instance, for the particular network trained in Figure 2b, the forecasted dynamic sea $\hat{\mu}_i$ level is within one standard deviation $\hat{\sigma}_i$ of the target value 62.7% of the time, and within $2\hat{\sigma}_i$ 90.8% of the time when evaluated over the test dataset. This is similar to how independent samples drawn from a true normal distribution will be within one standard deviation of the mean 68.3% of the time and within two standard deviations of the mean 95.4% of the time, indicating that forecasts are well-calibrated. Thus, the forecasted standard deviation gives an estimate of the acceptable level of forecast er-

ror. For this reason, we refer to predictions with a lower $\hat{\sigma}$ as "lower-uncertainty" or "higher confidence" predictions, although it will in general need to be verified that such low-uncertainty predictions result in lower errors on average.

Because the dynamic sea level forecasts are given as Gaussian distributions, these forecasts can be used to evaluate probabilities of specific events. For a given probabilistic event $A$, the implied probability by the neural network framework is

$$\hat{P}_i(A) = \hat{P}_i(Y \in A) = \int_A f(y|\hat{\mu}_i, \hat{\sigma}_i^2)\,dy \qquad (4)$$

where $Y$ is the random variable representing the target dynamic sea level and $f(y|\hat{\mu}_i, \hat{\sigma}_i^2)$ is the probability density function of a normal distribution with parameters $\hat{\mu}_i$ and $\hat{\sigma}_i^2$. One class of events that is useful to examine to identify sources of predictability is exceedance events. For instance, the implied probability of a positive sea level anomaly is given by

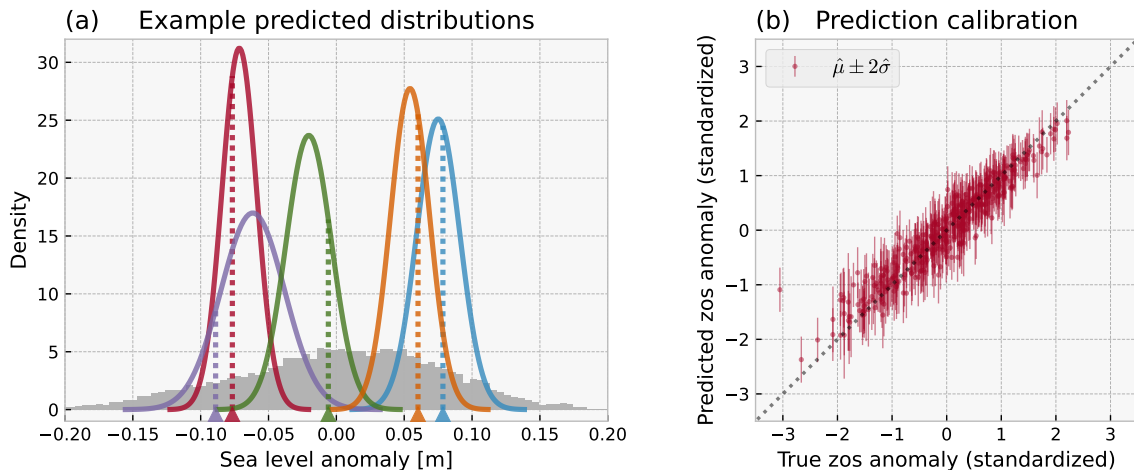$$\hat{P}_i(Y \geq 0) = \int_0^{\infty} f(y|\hat{\mu}_i, \hat{\sigma}_i^2)\,dy. \qquad (5)$$

FIG. 2: Example uncertainty-quantifying predictions given by a network trained at Guam (14°N, 145°E) for forecast lead $\tau = 20$ days. Panel (a) shows five example distributions predicted by the machine learning framework on the test set (solid lines), compared to the climatological distribution (gray histogram). The markers on the x-axis and vertical dotted lines indicate the true target sea level anomaly. Panel (b) shows a scatterplot of 200 random dynamic sea level predictions (mean ± two standard deviations) compared to the target dynamic sea level anomaly.

Nevertheless, whether the forecasted Gaussians can be used to skillfully predict probabilistic events must also be checked.

Neural networks trained on maximum-likelihood based loss functions were first proposed by Nix and Weigend (1994) and have been previously used several times in climate science to estimate aleatoric uncertainties (Guillaumin and Zanna 2021; Barnes and Barnes 2021; Gordon and Barnes 2022; Barnes et al. 2023). Typically, an assumption is made about the underlying sampling distribution for the uncertainties, and networks are trained to predict all parameters of the probability distribution together. In practice, training times were much shorter when training a network for the mean and standard deviation separately—an important practical consideration when training networks at several thousand locations to build spatial coverage. This two-stage training procedure is similar to the approach in Adler and Öktem (2018) and Perezhogin et al. (2023). We also note that using the MSE as a training loss for the point-estimate network, as opposed to other popular regression loss functions such as the Mean Absolute Error (MAE) or Huber loss (Huber 1964), is appropriate from a theoretical perspective, as the maximum-likelihood estimator for the mean under independent, identically distributed Gaussian samples is also the argument minimizer of the MSE. Section 3a explores the validity of the two-stage training procedure by assessing the networks' probabilistic predictions.

### c. Baselines: damped persistence and simple logistic regression

As a baseline, we compare our networks to forecasts produced by damped persistence (DP, Lorenz 1973). Given local observations of dynamic sea level at a given time $x(t) \in \mathbb{R}$, the damped persistence forecast at lead $\tau$, $\hat{x}(t + \tau) \in \mathbb{R}$, is given by

$$\hat{x}(t + \tau) = \beta_\tau x(t), \qquad (6)$$

where the autocorrelation coefficient $\beta_\tau \in [0, 1]$ is the best intermediate between a climatological forecast ($\beta_\tau = 0$) and a pure persistence forecast ($\beta_\tau = 1$). Thus, the damped persistence forecast encompasses the local linear predictability of sea level. Note that the features used in the damped persistence model have the same resolution as the target, so while the neural networks use input features at 5° resolution to predict sea level at 1°, the damped persistence model uses 1° inputs and outputs. The autocorrelation coefficient is computed via least squares estimation to minimize $\|x(t_i + \tau) - \beta_\tau x(t_i)\|_2$ over all times $t_i$. Autocorrelation coefficients are computed for each ensemble member in the training set and then averaged over all members.

The damped persistence model is simple and is commonly used for demonstrating deterministic forecast skill. However, because our uncertainty-quantifying framework is probabilistic, we also consider a probabilistic baseline to assess probabilistic forecasts. We employ a logistic regression baseline which uses the same features as the damped persistence model. Given an event $A$, the forecasted prob-

ability of the event occurring $\hat{P}_A(t+\tau)$ is given by

$$\hat{P}_A(t+\tau) = \frac{1}{1 + e^{-\left(\beta_\tau^{(0)} + \beta_\tau^{(1)} x(t)\right)}} \tag{7}$$

where $\beta_\tau^{(0)}$ and $\beta_\tau^{(1)}$ are coefficients also determined by least-squares. The logistic regression baseline, like the damped persistence model, characterizes local predictability of dynamic sea level, but for probabilistic events. The coefficients of the logistic regression model are fitted using the Scikit-learn package (Pedregosa et al. 2011) with the default $l_2$ regularization weight 1.

## 3. Results

### a. Forecast performance

1) DETERMINISTIC PERFORMANCE AND STATE-DEPENDENT PREDICTABILITY

Metrics evaluating the deterministic performance of all 6,590 of the networks trained at forecast leads of $\tau = 20$ and $\tau = 120$ days are shown in Figure 3. At forecast leads of $\tau = 20$ days (Fig. 3a), the lowest mean absolute errors (MAE) occur primarily in the low-latitude Pacific and Indian Oceans, and the Pacific and Indian Ocean eastern boundaries bordering North America and Australia. ANN errors are high in the Southern Ocean, where ocean dynamics are dominated by baroclinic instabilities. Low forecast errors persist up to leads of $\tau = 120$ days in the eastern and western tropical Pacific and tropical Indian Ocean.

To contextualize the ANN performance, Figures 3c and 3d compare the MAE of the ANN predictions to the errors of the damped persistence baseline described in Section 2c. The plots show the difference between damped persistence and ANN mean absolute errors at each location, $\Delta\text{MAE} = \text{MAE}_{DP} - \text{MAE}_{ANN}$, so that positive values indicate that the ANN has skill relative to damped persistence. Even by forecast leads of $\tau = 20$ days, the neural networks outperform damped persistence forecasts in the majority (91.7%) of locations (Fig. 3c). Regions where the ANN does not outperform damped persistence mostly occur along western boundary currents like the Gulf Stream, Kuroshio and Agulhas Current, where spatial gradients are relatively strong and the coarse-resolution inputs of the ANN may not accurately describe the local dynamic sea level. Nevertheless, since damped persistence forecasts represent the local, linear predictability of sea level, the prevalence of positive $\Delta\text{MAE}$ indicates that nonlocal and/or nonlinear dynamics impact dynamic sea level on daily-to-seasonal timescales in most locations. The most prominent regions of high ANN performance occur adjacent to the equator, possibly due to equatorial Rossby or Kelvin waves, which propagate dynamic sea level anomalies parallel to the equator. As forecast leads increase to $\tau = 120$ days, the regions of high skill expand towards higher latitudes throughout the tropics. This could reflect the fact that the phase velocity of Rossby waves decreases with latitude (Salmon 1998; Rossby 1939), so that longer lead times are needed to transmit nonlocal sea level signals to the forecast location. The proportion of locations where the ANN outperforms damped persistence also increases to 98.4%.

As mentioned previously, the purpose of the uncertainty-quantifying network is to identify initial conditions which may result in better forecasts. However, whether the networks are truly able to identify such initial conditions must be verified. Therefore, Figures 3e and 3f show the difference between ANN MAE computed over all predictions and the MAE of the predictions of only the 20% most-confident predictions from the test dataset. In 99.4% of locations at leads of $\tau = 20$ days and 98.2% of locations at leads of $\tau = 120$ days, the MAE for the 20% of predictions deemed the most-confident by the uncertainty-quantifying network is lower than the MAE over all predictions. Thus, the networks are able to identify state-dependent predictability at nearly all locations on these timescales. Notably, many of the regions of the largest state-dependent predictability identified by the uncertainty-quantifying network occur in complementary regions to the ANN average skill relative to damped persistence, such as in the midlatitude Pacific and Southern Ocean. Therefore, the most confident predictions made by the neural network framework are often significantly better than the average damped persistence prediction.

Figure 4 clarifies the relationship between predicted uncertainty and prediction errors for networks trained at Guam (14°N, 145°E) and in the western Indian Ocean (11°S, 60°E), indicated by the green dots labelled "A" and "B," respectively, in Figure 3. While the correspondence is not exact, the predicted uncertainty approximates the prediction MAE for each confidence decile. Figure 4 also demonstrates the utility of focusing on state-dependent predictability. For instance, although the average point-estimate neural network prediction is marginally better than damped persistence at forecast leads of $\tau = 20$ days at both locations, the most-confident predictions are significantly better. As forecast leads increase to $\tau = 120$ days, the difference between the neural network framework and damped persistence is even more apparent. Not only does the average network prediction outperform the baseline, but also the difference between the errors for the most-confident predictions and the median-confidence predictions increases for both of these locations.

2) PROBABILISTIC PREDICTIONS

To complement the deterministic evaluations of the point-estimate network in the previous section, here we evaluate the networks' prediction performance in a probabilistic sense. To assess the probabilistic forecasts,
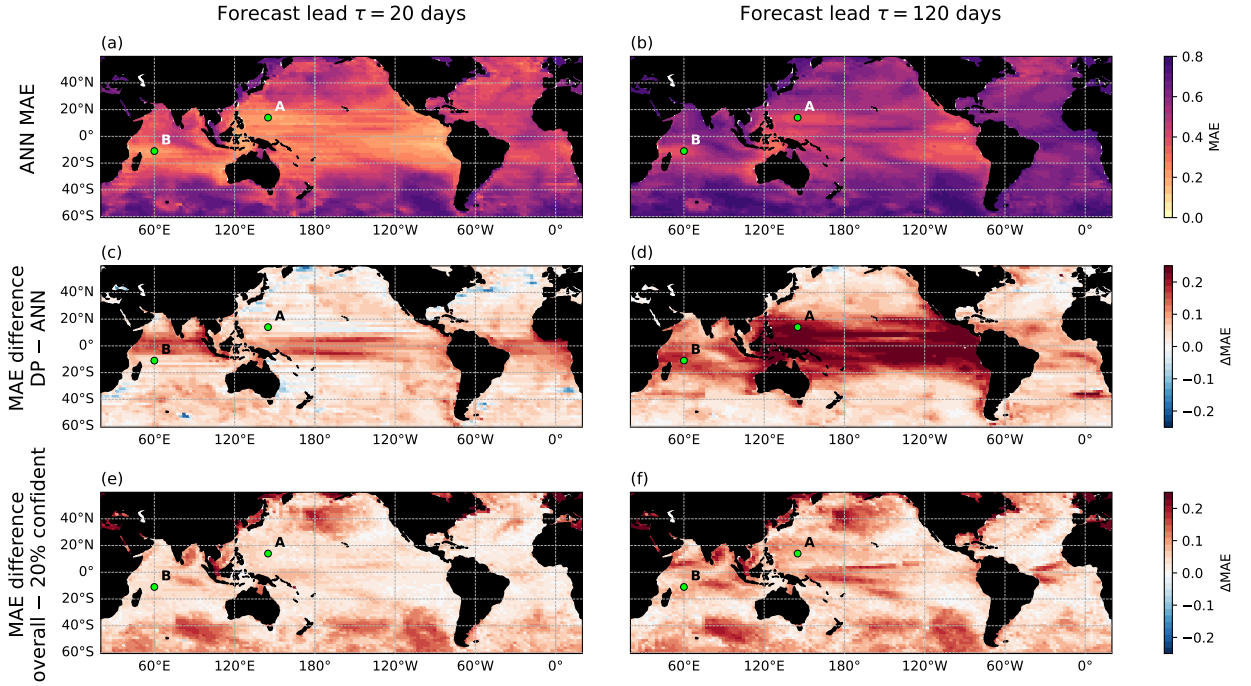
FIG. 3: Global ANN prediction error metrics for forecast leads of $\tau = 20$ days (a, c, e) and $\tau = 120$ days (b, d, f). Note that predictions are made at every other gridpoint latitude and longitude on the nominal $1°$ grid, but visualized here as a continuous map using nearest-neighbor interpolation. (a, b) Standardized mean absolute errors of the point-estimate networks over all samples. (c, d) Difference between damped persistence MAE and ANN MAE (positive values indicate ANN outperforms damped persistence. (e, f) Difference in ANN mean absolute errors taken over all samples and mean absolute errors taken over the 20% most-confident predictions as decided by the uncertainty-quantifying network.

we compute the Continuous Ranked Probability Score (CRPS) of the forecasts averaged over all of the samples (Matheson and Winkler 1976; Gneiting and Raftery 2007; Bröcker 2012). Given a set of predicted cumulative distribution functions $\hat{\boldsymbol{F}} = (\hat{F}_1, \ldots, \hat{F}_N)$ and observations $\boldsymbol{y} = (y_1, \ldots, y_N)$, the CRPS is defined by the $L^2$-distance between the predicted and empirical distributions:

$$\text{CRPS}(\hat{\boldsymbol{F}}, \boldsymbol{y}) = \frac{1}{N} \sum_{i=1}^{N} \int_{-\infty}^{\infty} \left( \hat{F}_i(x) - \mathbb{1}_{\{x > y_i\}} \right)^2 \, dx. \quad (8)$$

The predicted cumulative distribution functions given by the neural network are Gaussian with parameters determined by the neural networks $\hat{F}_i(y) = \Phi(\hat{\sigma}_i^{-1}(y - \hat{\mu}_i))$, where $\Phi$ is the cumulative distribution function of the standard normal distribution.

The CRPS is often used for assessing probabilistic forecasts against observations. It is a strictly-proper scoring rule, meaning that probabilistic forecasts cannot artificially improve CRPS by hedging for different outcomes. Thus, CRPS penalizes overconfident predictions as well as underconfident predictions. To normalize CRPS at each location, we compute a continuous ranked probability skill

score (CRPSS) relative to climatology by

$$\text{CRPSS}(\hat{\boldsymbol{F}}, \hat{\Psi}) = 1 - \frac{\text{CRPS}(\hat{\boldsymbol{F}}, \boldsymbol{y})}{\text{CRPS}(\hat{\Psi}, \boldsymbol{y})}, \quad (9)$$

where $\hat{\Psi}$ is the probabilistic forecast given by assuming a normal distribution with the climatological mean and standard deviation. Thus, CRPSS close to 1 indicates perfect forecast skill, while CRPSS close to 0 indicates predictions no better than climatological forecasts. The CRPS is computed using the properscoring module in Python 3 (Barrett et al. 2015).

CRPSS for all networks trained at forecast leads of $\tau = 20$ days and $\tau = 120$ days are shown in Figure 5a and 5b, respectively. Regions of high CRPSS occur mostly in the low-latitude Pacific and Indian Ocean. The spatial distribution of CRPSS is very highly correlated to the maps of MAE in Figure 3a and 3b ($R^2 = 0.988$ at leads $\tau = 20$ days and $R^2 = 0.963$ at leads $\tau = 120$). Such high correlations are expected, due to the fact that the MAE is a special case of the CRPS for deterministic predictions, where $\hat{F}_i(y)$ is represented by Heaviside functions $H(y - \hat{y}_i)$.
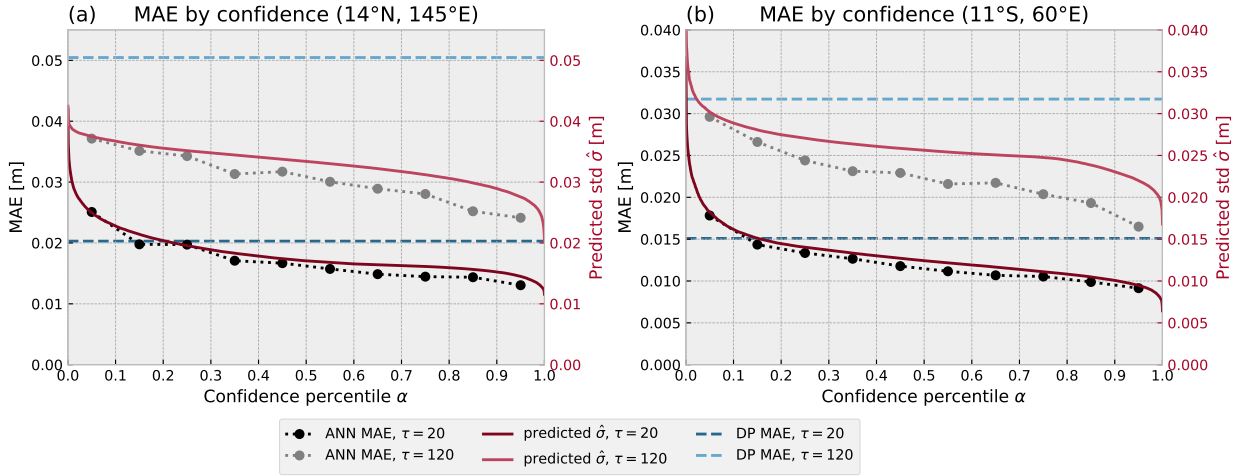
FIG. 4: Test dataset mean absolute error by network confidence level for networks trained at Guam (14°N, 145°E, panel a) and in the western Indian Ocean (11°S, 60°E, panel b). The red curve shows the predicted standard deviation outputted by the uncertainty-quantifying network as a function of the network confidence level (which is defined implicitly by the percentiles of predicted standard deviations in the test set). The black circles show the network MAEs for predictions grouped by network confidence decile. The blue line shows the MAE for the damped-persistence model. Dark colors indicate errors and uncertainties for forecast leads of $\tau = 20$ days, whereas lighter shades indicate forecast leads of $\tau = 120$ days.



FIG. 5: Global probabilistic performance metrics for networks trained at forecast leads of $\tau = 20$ days (a, c) and $\tau = 120$ days (b, d). (a, b) Continuous ranked probability skill score (Eq. 9) of the neural network framework evaluated relative to climatology. (c, d) Brier skill scores (Eq. 11) for predictions of positive anomalies made by the networks evaluated relative to the logistic regression baseline discussed in Section 2c.

We also compute the Brier score, which is useful for assessing a model's ability to predict probabilistic events with binary outcomes (Brier 1950). The Brier score of a set of Bernoulli probability forecasts $\hat{\boldsymbol{P}} = (\hat{P}_1, \ldots, \hat{P}_N)$ against binary observations of class occurrences $\boldsymbol{o} =$

$(o_1, \ldots, o_N) \in \{0, 1\}$ is given by

$$\text{BS}(\hat{\boldsymbol{P}}, \boldsymbol{o}) = \frac{1}{N} \sum_{i=1}^{N} \left( \hat{P}_i - o_i \right)^2 \qquad (10)$$

A Brier skill score, BSS, relative to a baseline model $\hat{Q}$, may also be defined by

$$\text{BSS}(\hat{P}, \hat{Q}) = 1 - \frac{\text{BS}(\hat{P}, o)}{\text{BS}(\hat{Q}, o)} \qquad (11)$$

Like the CRPSS, BSS near 1 implies perfect probabilistic forecasts, while BSS near 0 indicates forecasts no better than the baseline.

Figure 5c and 5d shows the Brier skill scores of the neural networks' ability to predict positive anomalies (as defined by Equation 5) relative to the logistic regression baseline. Brier skill scores are greater than 0 at most locations (65.8% of locations at leads of $\tau = 20$ days and 70.0% of locations at $\tau = 120$ days). That is, although the neural network framework has not been explicitly tasked with predicting positive sea level anomaly events, probabilities implied using the networks' predicted mean and standard deviation are better at predicting positive sea level anomalies than a logistic regression baseline which has been explicitly trained for this purpose at most locations. Of course, the neural network uses spatial information which is not available to the logistic regression model. Nevertheless, the skill of the neural networks for predicting positive anomalies illustrates how the uncertainty-quantifying neural network framework can be used for predicting exceedance events. Regions of high BSS expand from the low-latitude Indo-Pacific at forecast leads of $\tau = 20$ days towards higher latitudes at $\tau = 120$. This indicates regions in which nonlocal factors become important for forecasting exceedance probabilities.

Due to the high CRPSS and BSS and low MAE at the locations labeled "A" and "B" in Figure 3 and 5 (14°N, 145°E and 11°S, 60°E, respectively), the remainder of this paper focuses on these locations. Particular attention is given to identifying how drivers of sea level predictability change over different daily-to-seasonal forecast leads and what these primary drivers are.

### b. Predictability by forecast lead

Figure 6 shows the performance of each of the deterministic forecasting techniques on a variety of lead times from $\tau = 10$ to $\tau = 180$ days at Guam (location "A", 14°N, 145°E) and the western Indian Ocean (location "B", 11°S, 60°E). Over daily-to-seasonal timescales, damped persistence errors decay towards climatological errors at both locations. Although the predictions for the neural network have similar errors to damped persistence at forecast leads of 10 or 20 days, ANN errors grow more slowly. ANN skill peaks at about 120 days at Guam with 40.2% improvement in MAE over damped persistence, while it peaks at 28.5% skill at 60 days in the western Indian Ocean. The 20% most-confident predictions also have lower MAE than the

average predictions at all time lags. This shows that focusing on state-dependent predictability can consistently offer advantages for forecasts at these locations.

The MAE of forecasts in Figure 6 increase with forecast lead because of the degradation of initial condition information from boundary forcing and chaotic dynamics. Thus, the average predicted uncertainties can be expected to also increase with larger forecast leads, plateauing at the climatological uncertainty when initial condition information has fully deteriorated. Figure 7 shows the distribution of predicted uncertainties by the neural network as a function of the forecast lead. The average predicted standard deviation does indeed increase with forecast lead, closely matching the mean absolute errors of the point-estimate network predictions and remaining less than the climatological standard deviation at all forecasting leads. Moreover, not only do the *mean* absolute errors increase with forecast leads but also the *spread* of the distribution of absolute errors increases with forecast leads. Accordingly, the range of predicted standard deviations increases with forecast lead. Using a $t$-distributed Wald test for positive slope finds that the increasing minimum-maximum and 90%-interpercentile ranges of predicted standard deviations are significant at the 5% significance level for both locations.

The increasing MAE and broadening range of errors due to the deterioration of initial condition information also have implications for the predicted probabilities of exceedance events. For instance, predictions with completely certain exceedance outcomes would be represented by exceedance probabilities of exactly 0 or exactly 1. On the other hand, forecasts of exceedance probabilities in which initial conditions provide no information about the outcome would simply output climatological probabilities of the event. The amount of information contained in the initial conditions for each forecast is conveyed by the empirical distributions of predicted probabilities for positive anomalies in Figure 8. At shorter forecast leads of $\tau = 20$ days, most of the neural network predicted probabilities of positive sea level anomalies are near 0 or 1: 68% of ANN predicted probabilities $\hat{P}_i$ yield $\hat{P}_i < 0.05$ or $\hat{P}_i > 0.95$ for location A and 63% of probabilities are this extreme for location B. However, as forecast leads increase, the proportion of high-confidence predictions of exceedance probabilities decreases: at leads of $\tau = 120$ days, this proportion of extreme probabilities has fallen to 47% for location A and 30% for location B. The loss of information contained in the initial conditions for the logistic regression model is much more pronounced. For instance, while the number of highly certain predictions at location A yielding probabilities of less than 0.05 or greater than 0.95 of the logistic regression model is similar to the ANN at leads of $\tau = 20$ days (56% for logistic regression vs 67% for the ANN), by forecast leads of $\tau = 120$ days, only 1% percent of predictions are this confident for the logistic regression model
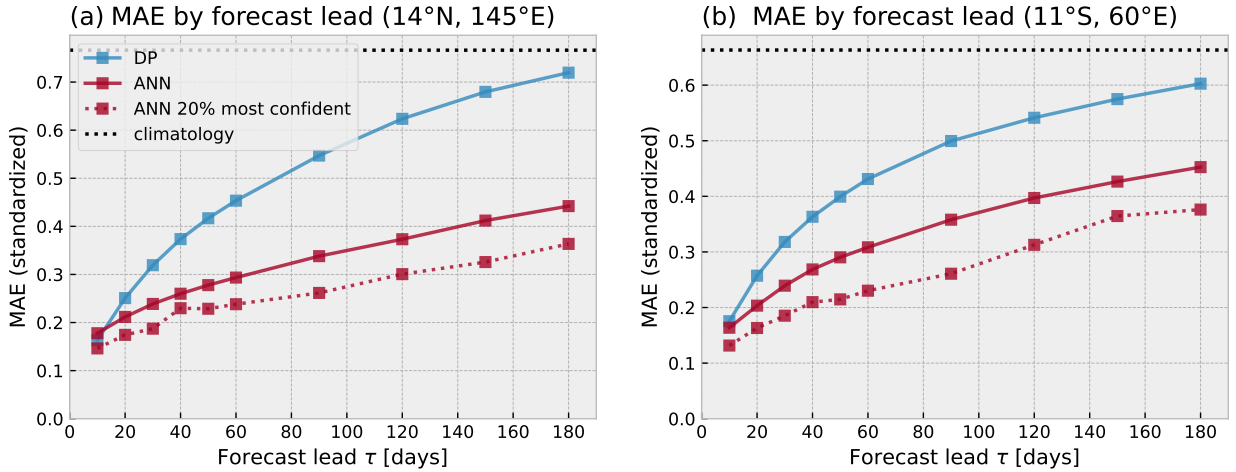
FIG. 6: Mean absolute error by forecast lead time for networks trained at Guam (14°N, 145°E, panel a) and in the western Indian Ocean, (pointB, panel b). Markers indicate forecasting leads at which networks have been trained. Black dashed line shows the climatological error, while the blue curve shows the damped persistence error as a function of lead time. The red solid line shows the Mean Absolute Error of the ANN over all samples, and the red dashed line shows the errors of only the 20% most confident predictions.
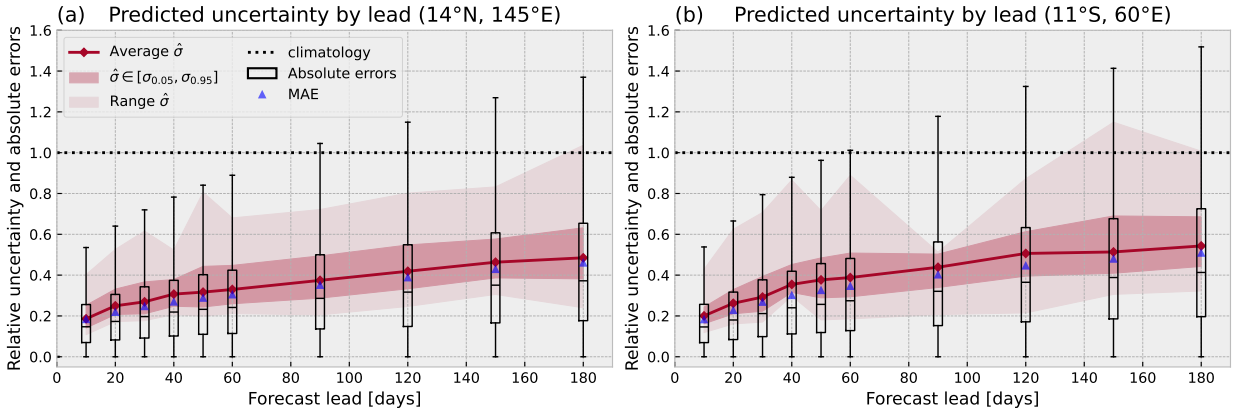


FIG. 7: Distribution of predicted standard deviations $\hat{\sigma}_i$ as a function of forecast lead at Guam (14°N, 145°E, panel a) and in the Western Indian Ocean (11°S, 60°E, panel b). Predicted uncertainties are normalized relative to the climatological standard deviation. Solid red line indicates the average predicted uncertainty over all samples in the test set for each forecasting lead. Dark red shading indicates the range of standard deviations between the 5th and 95th percentiles, whereas light red shading indicates the entire min-max range of predicted standard deviations. The boxplots show the distribution of the absolute errors between the true and forecasted sea level using the point-estimate network with mean absolute errors indicated using blue triangles (outliers are removed for visualization purposes).

whereas 47% percent of the ANN predictions are this confident. The fact that initial condition information decays more quickly for the logistic regression baseline than for the neural networks suggests that non-local information becomes more important for predicting exceedance probabilities while local information becomes less important over daily-to-seasonal timescales.

*c. Drivers of predictability over different forecast leads*

Figure 8 suggests that both local and nonlocal drivers can impact sea level predictability on daily-to-seasonal timescales. To identify these drivers, Figures 9 and 10 show the composite inputs averaged over the 20% of samples most likely to result in positive sea level anomalies as predicted by the neural networks at Guam (14°N, 145°E)
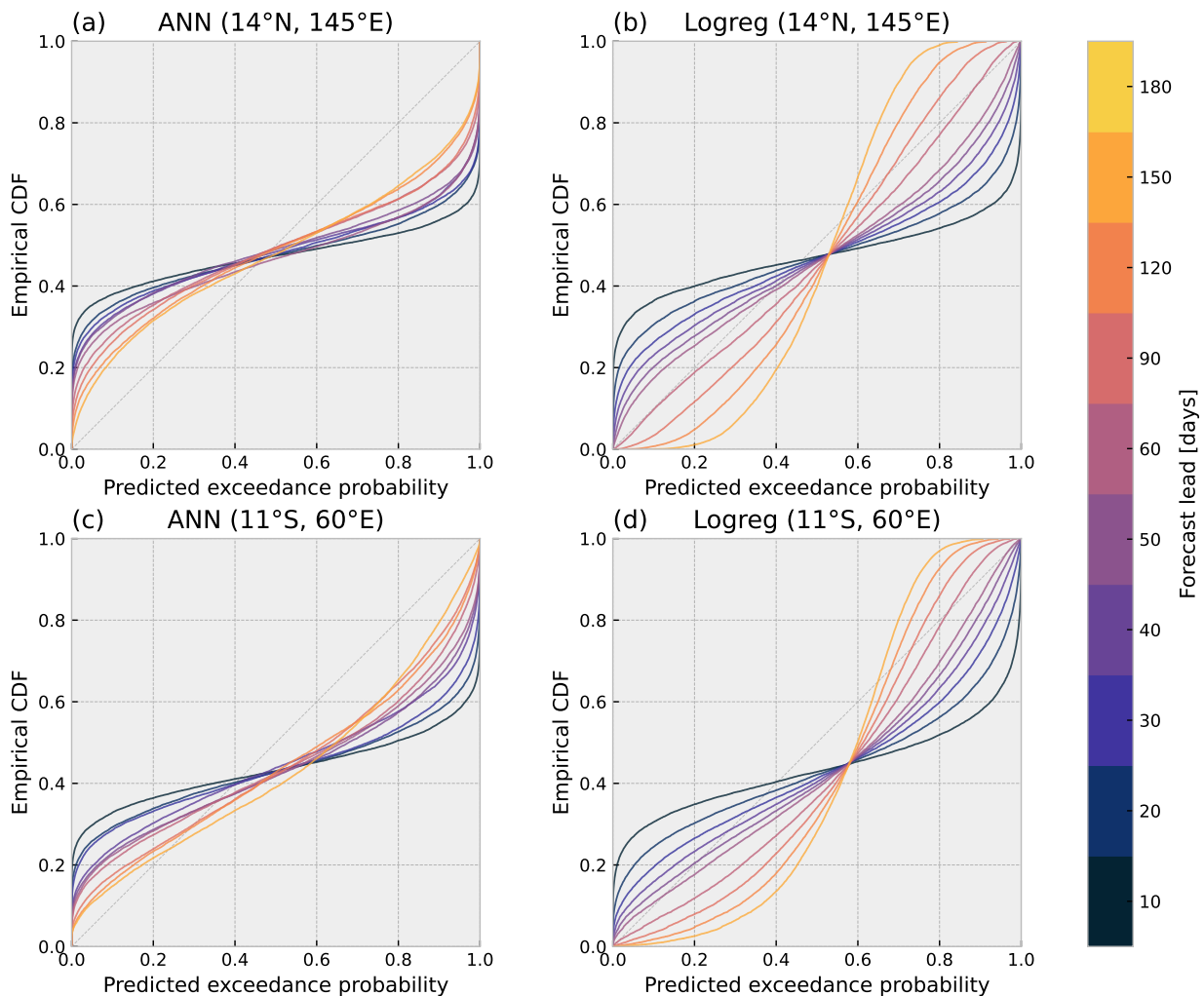
FIG. 8: Empirical distribution of predicted probabilities of positive sea level anomaly exceedance events for neural network framework (a, c) and logistic regression baseline (b, d) at Guam, 14°N, 145°E (a, b) and in the Western Indian Ocean 11°S, 60°E (c, d).

and the western Indian Ocean (11°S, 60°E), respectively. (While thresholding on exceedance probability is also possible, it is avoided as it yields an inconsistent number of samples for different time lags.) The composites thus show the average initial condition resulting in a likely positive anomaly (using the definition of exceedance probabilities from Equation 5) . While the composites show *which* samples result in likely positive anomalies, it is unclear *why* the networks have selected these samples. Therefore, a neural network attribution technique, integrated gradients (Sundararajan et al. 2017; Mamalakis et al. 2022), was applied to the point-estimate networks to quantify the relative contribution of each input feature to positive sea level anomalies. Integrated gradients are applied to the point-estimate network to identify input features that result in

more positive anomalies. The 95[th]-percentile integrated gradients for each variable and forecast lead are stippled in Figures 9 and 10.

At Guam (14°N, 145°E, Fig. 9), the most prominent initial conditions resulting in likely positive anomalies is the persistence of local sea level anomalies at all time lags. The local dynamic sea level pattern resulting in likely positive predictions is stretched significantly zonally, extending especially eastward from the prediction location. Moreover, the regional maximum-composite input dynamic sea level moves eastward with increasing forecast lead. This could be due to incident Rossby waves, which propagate sea level signals with a westward phase velocity and could be a non-local source of predictability (Vallis 2017).
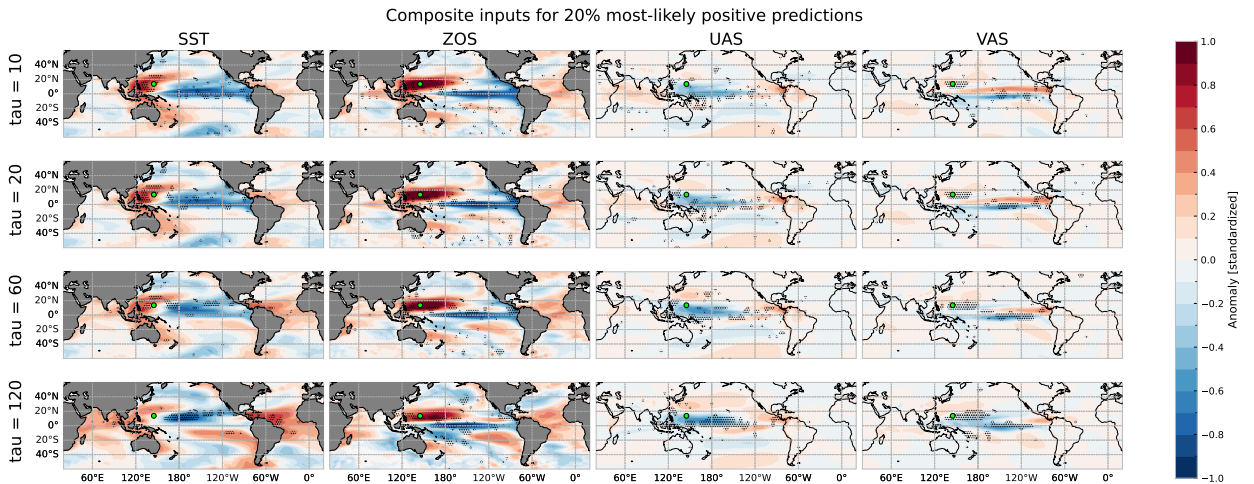
FIG. 9: Composites averaged over the 20% of input samples most likely to result in positive sea level anomalies at Guam (14°N, 145°E, green dot) for different forecast leads. Columns show different input fields (SST, ZOS, UAS, VAS), and rows show different forecast leads ($\tau$ = 10, 20, 60, and 120 days). Stippling indicates integrated gradients of the point-estimate network exceeding the 95[th] percentile for each variable.

The SST fields resulting in likely-positive sea level anomalies at Guam exhibit a pattern that persists over forecast leads from 10–120 days, with positive SST anomalies northwest of the prediction location but with negative and intensifying SSTs eastward in the Central Pacific. This SST pattern is somewhat distinct from some major drivers of climate variability occurring in the Pacific, lacking the signature east Pacific tongue of the El Niño-Southern Oscillation or the strong footprint of SSTs in the central North Pacific of the Pacific Decadal Oscillation. However, the SST pattern is quite similar to the region of strong seasonal lagged-correlation between Pacific SSTs and sea level recorded by tide gauges at Guam found in Chowdhury et al. (2007b). Thus, larger surface ocean heat content in the western tropical Pacific drives positive sea level anomalies at Guam through dynamical mechanisms robust to both CESM2 simulations and observations.

In the western Indian Ocean (11°S, 60°E, Fig. 10), SST and ZOS composites point to a strong signal from the Indian Ocean resembling the positive phase of the Indian Ocean Dipole (IOD, Saji et al. 1999; Webster et al. 1999). The surface wind divergence east of Indonesia found in the composites is consistent with the Walker circulation that accompanies the dipole SST patterns. This indicates that the IOD is a significant source of information for the likelihood of sea level anomalies in the western Indian Ocean. Increases in SSTs in the western Indian Ocean during a positive IOD likely drive thermosteric sea level, driving persistent sea level anomalies and increasing the forecasted likelihood of positive anomalies on daily-to-seasonal timescales. The relationship is corroborated by Roberts et al. (2016), which identifies a significant finger-

print of the IOD on observed thermosteric sea level in the Indian Ocean. The dipole signature in SST and ZOS also slightly diminishes over the daily-to-seasonal timescale. This is consistent with the timescale of the IOD index, which ranges from weeks to months (Wang et al. 2016; Behera et al. 2013; Rao and Yamagata 2004).

El Niño also emerges as a prominent source of dynamic sea level predictability on daily-to-seasonal timescales, as evidenced by the intensifying central and eastern Pacific SST pattern. The integrated gradients indicate that while local SSTs are more important for predicting dynamic sea level at forecast leads of $\tau$ = 10 and $\tau$ = 20 days, tropical SSTs in the Niño 3.4 region are more useful for forecasting likely positive sea level anomalies at $\tau$ = 60 and $\tau$ = 120 days. El Niño events may impact Indian Ocean dynamic sea level through its basin-scale influence (Xie et al. 2002) as well as by its own influence on the IOD (Stuecker et al. 2017; Yang et al. 2015; Krishnamurthy and Kirtman 2003). Thus, even though the local predictability of western Indian Ocean dynamic sea level due to the IOD may decrease with forecast leads, El Niño may provide a remote source of predictability for longer leads by being a precursor to the IOD.

## 4. Discussion

Dynamic sea level on daily-to-seasonal timescales is driven by a variety of processes in the atmosphere and ocean. Identifying conditions where predictability is enhanced can reduce uncertainties and result in "windows of opportunity" for better forecasts. In this study, uncertainty-quantifying regression neural networks were trained on

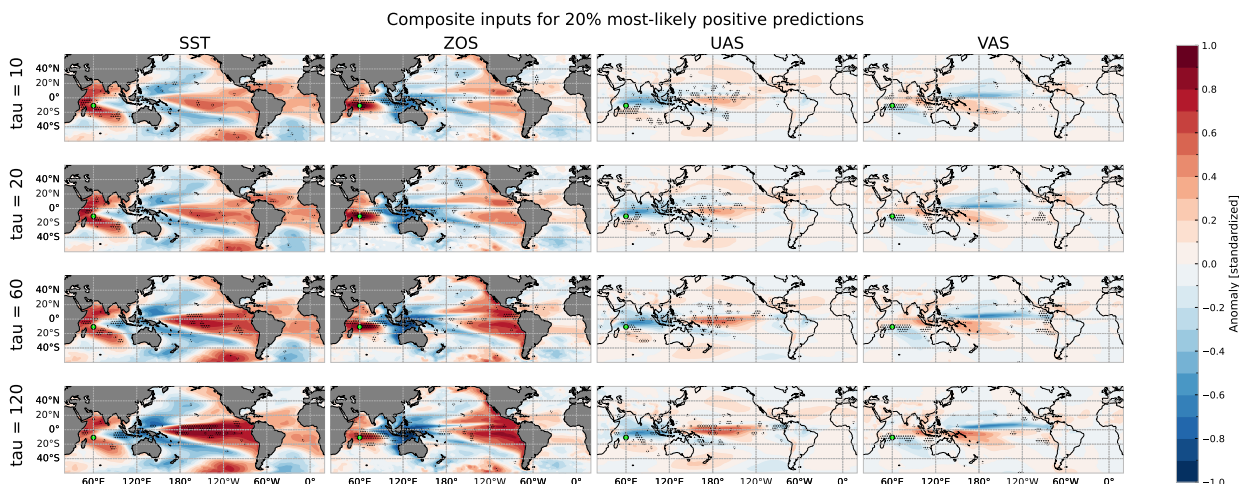Composite inputs for 20% most-likely positive predictions



Fig. 10: Same as Figure 9, but for predictions in the western Indian Ocean (11°S, 60°E).

CESM2 coupled climate data using the Gaussian maximum log-likelihood to identify sources of state-dependent predictability for dynamic sea level. The uncertainty-quantifying networks' predicted standard deviations can be leveraged to identify initial conditions that result in better forecasts at most locations and timescales. This establishes that state-dependent sources of predictability of sea level exist and demonstrates the utility of detecting these fortuitous initial conditions.

Although the networks discriminate different uncertainty levels based on initial states, the range of predicted standard deviations is arguably rather small. For example, 90% of the predicted standard deviations for the network trained at Guam at lead $\tau = 120$ days fall between 33% and 55% of the climatological standard deviation (Fig. 7). While the distribution of absolute errors may be statistically consistent with these predicted standard deviations, it is possible that using different resolutions, input features, or architectures could result in broader ranges of predicted standard deviations. Having a larger spread of predicted standard deviations could result in better estimates of the errors and enhance the utility of focusing on state-dependent sources of predictability.

In the uncertainty-quantifying framework, it was assumed that the uncertainties could be described using a Gaussian distribution and that the point-estimate and uncertainty-quantifying networks could be trained separately instead of learning all the parameters of the aleatoric distribution simultaneously. The relatively high CRPSS in the low-latitude Indo-Pacific suggests that these assumptions are adequate in these regions. The regions of low CRPSS occurring in the midlatitudes and the Southern Ocean do not necessarily mean that these choices are invalid. Indeed, the extremely high correlation between CRPSS and MAE imply a general lack of predictability

in these regions, resulting in systematic underprediction for the point estimate network (Murphy 1973) and overall predicted distributions similar to climatological forecasts.

Although not explicitly trained to forecast probabilities, we noted that the uncertainty-quantifying network framework implies forecasted exceedance probabilities that often outperform local logistic regression baselines, which *are* explicitly constructed to predict exceedance probabilities. This reflects not only the relative quality of the probabilistic predictions and validity of assuming normal-distributed uncertainties but also the importance of nonlocal sources of predictability, which are registered in the features of the neural networks but not represented in the inputs of the damped persistence or logistic regression baselines.

Using a regression neural network to forecast event probabilities also illustrates a useful application of the regression neural network framework of Gordon and Barnes (2022). Mayer and Barnes (2021) and Gordon and Barnes (2022) used different criteria to isolate sources of state-dependent predictability: whereas Gordon and Barnes (2022) used predicted standard deviation to define predictable initial conditions, Mayer and Barnes (2021) focused on discrete probabilities such as exceedance events. The different architectures of the two studies make these different notions of predictability natural, as the regression networks of Gordon and Barnes (2022) output standard deviation while the classification networks of Mayer and Barnes (2021) output event probabilities. However, as shown in this study, the regression framework of Gordon and Barnes (2022) can be used to predict event probabilities, potentially making it a more flexible approach to identifying state-dependent sources of predictability. A direct comparison of the skill of the different approaches for predicting event probabilities could clarify the strengths and limitations of the regression framework.

The distribution of predicted exceedance probabilities by the neural networks and logistic regression baseline conveys the degradation of information available from the initial conditions for predicting exceedance events. For both forecast techniques, near-certain probabilistic exceedance outcomes decrease in frequency as the forecast leads increase, indicating timescales where initial condition information is lost. However, the empirical distributions devolve to climatology more slowly for the neural networks than for the logistic regression, suggesting that regional-scale information about the predictability of positive sea level anomalies continues to persist while the local information fades. This approach may be a useful way to characterize the loss of predictability on daily-to-seasonal timescales. However, it should be stressed that the predicted probabilities depend not only on the initial conditions but also on the quality of the forecast models themselves.

Potential physical sources of dynamic sea level predictability were investigated by analyzing composites of the input samples deemed most likely by the neural networks to result in positive sea level anomalies. At Guam, propagating Rossby waves were identified as a potential source of predictability for sea level, while in the western Indian Ocean, the persistence of sea level anomalies due to the Indian Ocean Dipole provided a source of predictability, though the influence of El Niño begins to emerge on seasonal timescales.

Many of the sources of state-dependent predictability identified in the composites seem to come from low-frequency drivers, such as sea surface temperatures or dynamic sea level. Significant sources of predictability from the surface wind fields were more difficult to identify through our approach. Multiple studies have shown how surface winds can impact dynamic sea level through, for instance, manometric changes from wind stress (Hermans et al. 2022; Arcodia et al. 2024; Fukumori et al. 1998), surface Ekman mass convergence (Kamp et al. 2024; Piecuch and Ponte 2011), or changes in large-scale Sverdrup balance (Cabanes et al. 2006; Roberts et al. 2016). The prominence of low-frequency drivers found using our methods could be due to the averaging over multiple samples to identify robust initial states with predictable outcomes, as distinct, high frequency inputs would be filtered out. Cluster analysis methods, such as $k$-means or DBSCAN, could be used to identify distinct drivers. However, applying such methods introduces numerous sensitivities (e.g., type of clustering algorithm, dimensionality reduction techniques, and number of components) and is beyond the scope of this work. Nevertheless, the sources of state-dependent predictability for sea level identified in this study may help establish avenues for improving future sea level forecasts on daily-to-seasonal timescales.

## References

Adler, J., and O. Öktem, 2018: Deep bayesian inversion. *arXiv preprint arXiv:1811.05910*.

Albers, J. R., and M. Newman, 2019: A Priori Identification of Skillful Extratropical Subseasonal Forecasts. *Geophys. Res. Lett.*, **46 (21)**, 12 527–12 536.

Amaya, D. J., M. G. Jacox, J. Dias, M. A. Alexander, K. B. Karnauskas, J. D. Scott, and M. Gehne, 2022: Subseasonal-to-Seasonal Forecast Skill in the California Current System and Its Connection to Coastal Kelvin Waves. *J. Geophys. Res. Oceans*, **127 (1)**, e2021JC017 892.

Aparna, S., J. McCreary, D. Shankar, and P. Vinayachandran, 2012: Signatures of Indian Ocean Dipole and El Niño–Southern Oscillation events in sea level variations in the Bay of Bengal. *J. Geophys. Res. Oceans*, **117 (C10)**.

Arcodia, M. C., E. Becker, and B. P. Kirtman, 2024: Subseasonal Variability of US Coastal Sea Level from MJO and ENSO Teleconnection Interference. *Wea. Forecasting*, **39 (2)**, 441–458.

Balmaseda, M. A., R. McAdam, S. Masina, M. Mayer, R. Senan, E. de Bosisséson, and S. Gualdi, 2024: Skill assessment of seasonal forecasts of ocean variables. *Front. Mar. Sci.*, **11**, 1380 545.

Barnes, E. A., and R. J. Barnes, 2021: Controlled Abstention Neural Networks for Identifying Skillful Predictions for Regression Problems. *J. Adv. Model. Earth Syst.*, **13 (12)**, e2021MS002 575.

Barnes, E. A., R. J. Barnes, and M. DeMaria, 2023: Sinh-arcsinh-normal distributions to add uncertainty to neural network regression tasks: Applications to tropical cyclone intensity forecasts. *Env. Data Sci.*, **2**, e15.

Barrett, L., S. Hoyer, A. Kleeman, D. O'Kane, and Coauthors, 2015: properscoring. The Climate Corporation, gitHub, accessed 2024-12-12, https://github.com/properscoring/properscoring/tree/master.

Becker, M., B. Meyssignac, C. Letetrel, W. Llovel, A. Cazenave, and T. Delcroix, 2012: Sea level variations at tropical Pacific islands since 1950. *Global Planet. Change*, **80**, 85–98.

Behera, S., P. Brandt, and G. Reverdin, 2013: The tropical ocean circulation and dynamics. *Ocean Circulation and Climate*, Vol. 103, Elsevier, 385–412.

Brettin, A., 2025: Code for Brettin, Zanna, and Barnes (2025): daily-to-seasonal dynamic sea level predictabilty (v1.0.0). Zenodo, https://doi.org/10.5281/zenodo.14873345.

Brier, G. W., 1950: Verification of forecasts expressed in terms of probability. *Mon. Wea. Rev.*, **78 (1)**, 1–3.

Bröcker, J., 2012: Evaluating raw ensembles with the continuous ranked probability score. *Quart. J. Roy. Meteor. Soc.*, **138 (667)**, 1611–1617.

Cabanes, C., T. Huck, and A. Colin de Verdière, 2006: Contributions of Wind Forcing and Surface Heating to Interannual Sea Level Variations in the Atlantic Ocean. *J. Phys. Oceanogr.*, **36 (9)**, 1739–1750.

Chen, L., J. Yang, and L. Wu, 2023: Topography Effects on the Seasonal Variability of Ocean Bottom Pressure in the North Pacific Ocean. *J. Phys. Oceanogr.*, **53 (3)**, 929–941.

Chowdhury, M. R., P.-S. Chu, and T. Schroeder, 2007a: ENSO and seasonal sea-level variability–a diagnostic discussion for the US-Affiliated Pacific Islands. *Theor. Appl. Climatol.*, **88**, 213–224.

Chowdhury, M. R., P.-S. Chu, T. Schroeder, and N. Colasacco, 2007b: Seasonal sea-level forecasts by canonical correlation analysis—an operational scheme for the U.S.-affiliated Pacific Islands. *Int. J. Climatol.*, **27 (10)**, 1389–1402.

Christensen, H. M., J. Berner, and S. Yeager, 2020: The Value of Initialization on Decadal Timescales: State-Dependent Predictability in the CESM Decadal Prediction Large Ensemble. *J. Climate*, **33 (17)**, 7353–7370.

Computational and Information Systems Laboratory, 2023: Casper: HPE Cray EX System (University Community Computing). NSF National Center for Atmospheric Research, Boulder, CO, https://doi.org/10.5065/qx9a-pg09.

Danabasoglu, G., C. Deser, K. Rodgers, and A. Timmermann, 2021: CESM2 Large Ensemble Dataset. National Center for Atmospheric Research, URL https://www.earthsystemgrid.org/dataset/ucar.cgd.cesm2le.output.html, https://doi.org/https://doi.org/10.26024/kgmp-c556.

Danabasoglu, G., and Coauthors, 2020: The Community Earth System Model Version 2 (CESM2). *J. Adv. Model. Earth Syst.*, **12 (2)**, e2019MS001 916.

DeMott, C., Á. Muñoz, C. Roberts, C. Spillman, and F. Vitart, 2021: The benefits of better ocean weather forecasting. *Eos*, **102**.

Doi, T., M. Nonaka, and S. Behera, 2020: Skill Assessment of Seasonal-to-Interannual Prediction of Sea Level Anomaly in the North Pacific Based on the SINTEX-F Climate Model. *Front. Mar. Sci.*, **7**, 546 587.

Dukowicz, J. K., and R. D. Smith, 1994: Implicit free-surface method for the Bryan-Cox-Semtner ocean model. *J. Geophys. Res. Oceans*, **99 (C4)**, 7991–8014.

Eyring, V., S. Bony, G. A. Meehl, C. A. Senior, B. Stevens, R. J. Stouffer, and K. E. Taylor, 2016: Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6) experimental design and organization. *Geosci. Model Dev.*, **9 (5)**, 1937–1958.

Fasullo, J. T., P. R. Gent, and R. S. Nerem, 2020: Sea Level Rise in the CESM Large Ensemble: The Role of Individual Climate Forcings and Consequences for the Coming Decades. *J. Climate*, **33 (16)**, 6911–6927.

Fox-Kemper, B., and Coauthors, 2021: Climate Change 2021: The Physical Science Basis. Contribution of Working Group I to the Sixth Assessment Report of the Intergovernmental Panel on Climate Change. Tech. rep., Intergovernmental Panel on Climate Change, Cambridge, United Kingdom and New York, NY, USA.

Frame, T., J. Methven, S. Gray, and M. Ambaum, 2013: Flow-dependent predictability of the North Atlantic jet. *Geophys. Res. Lett.*, **40 (10)**, 2411–2416.

Fraser, R., M. Palmer, C. Roberts, C. Wilson, D. Copsey, and L. Zanna, 2019: Investigating the predictability of North Atlantic sea surface height. *Climate Dyn.*, **53**, 2175–2195.

Frederikse, T., and Coauthors, 2020: The causes of sea-level rise since 1900. *Nature*, **584 (7821)**, 393–397.

Fukumori, I., R. Raghunath, and L.-L. Fu, 1998: Nature of global large-scale sea level variability in relation to atmospheric forcing: A modeling study. *J. Geophys. Res. Oceans*, **103 (C3)**, 5493–5512.

Gill, A., and P. Niller, 1973: The theory of the seasonal variability in the ocean. *Deep-Sea Res. Oceanogr. Abstr.*, **20 (2)**, 141–177.

Gneiting, T., and A. E. Raftery, 2007: Strictly proper scoring rules, prediction, and estimation. *J. Am. Stat. Assoc.*, **102 (477)**, 359–378.

Gordon, E. M., and E. A. Barnes, 2022: Incorporating Uncertainty Into a Regression Neural Network Enables Identification of Decadal State-Dependent Predictability in CESM2. *Geophys. Res. Lett.*, **49 (15)**, e2022GL098 635.

Gregory, J. M., and Coauthors, 2019: Concepts and terminology for sea level: Mean, variability and change, both local and global. *Surv. Geophys*, **40**, 1251–1289.

Griffies, S. M., and Coauthors, 2014: An assessment of global and regional sea level for years 1993–2007 in a suite of interannual CORE-II simulations. *Ocean Model.*, **78**, 35–89.

Griffies, S. M., and Coauthors, 2016: OMIP contribution to CMIP6: Experimental and diagnostic protocol for the physical component of the Ocean Model Intercomparison Project. *Geosci. Model Dev.*, 3231.

Guillaumin, A. P., and L. Zanna, 2021: Stochastic-Deep Learning Parameterization of Ocean Momentum Forcing. *J. Adv. Model. Earth Syst.*, **13 (9)**, e2021MS002 534.

Hermans, T. H., C. A. Katsman, C. M. Camargo, G. G. Garner, R. E. Kopp, and A. B. Slangen, 2022: The Effect of Wind Stress on Seasonal Sea-Level Change on the Northwestern European Shelf. *J. Climate*, **35 (6)**, 1745–1759.

Hino, M., S. T. Belanger, C. B. Field, A. R. Davies, and K. J. Mach, 2019: High-tide flooding disrupts local economic activity. *Sci. Adv.*, **5 (2)**, eaau2736.

Hochet, A., W. Llovel, T. Huck, and F. Sévellec, 2024: Advection surface-flux balance controls the seasonal steric sea level amplitude. *Sci. Rep.*, **14 (1)**, 10 644.

Huber, P. J., 1964: Robust Estimation of a Location Parameter. *Ann. Math. Stat.*, **35 (1)**, 73–101.

Jacox, M. G., and Coauthors, 2020: Seasonal-to-interannual prediction of north american coastal marine ecosystems: Forecast methods, mechanisms of predictability, and priority developments. *Prog. Oceanogr.*, **183**, 102 307.

Kalnay, E., and A. Dalcher, 1987: Forecasting forecast skill. *Mon. Wea. Rev.*, **115 (2)**, 349–356.

Kamp, W., W. Han, L. Zhang, S. Kido, and J. P. McCreary, 2024: Tropical Atmospheric Intraseasonal Oscillations Leading to Sea Level Extremes in Coastal Indonesia during Recent Decades. *J. Climate*, **37 (9)**, 2867–2880.

Kenigson, J. S., W. Han, B. Rajagopalan, M. Jasinski, and Coauthors, 2018: Decadal Shift of NAO-Linked Interannual Sea Level Variability along the US Northeast Coast. *J. Climate*, **31 (13)**, 4981–4989.

Kingma, D. P., and J. Ba, 2014: Adam: A Method for Stochastic Optimization. *arXiv preprint arXiv:1412.6980*, https://doi.org/10.48550/arXiv.1412.6980.

Krishnamurthy, V., 2019: Predictability of weather and climate. *Earth Space Sci.*, **6 (7)**, 1043–1056.

Krishnamurthy, V., and B. P. Kirtman, 2003: Variability of the Indian Ocean: Relation to monsoon and ENSO. *Quart. J. Roy. Meteor. Soc.*, **129 (590)**, 1623–1646.

LeCun, Y., L. Bottou, G. B. Orr, and K.-R. Müller, 2002: Efficient BackProp. *Neural Networks: Tricks of the Trade*, G. M. et al., Ed., Springer, 9–48.

Li, S., and Coauthors, 2022: Contributions of Different Sea-Level Processes to High-Tide Flooding Along the U.S. Coastline. *J. Geophys. Res. Oceans*, **127 (7)**, e2021JC018 276.

Lin, S.-J., and R. B. Rood, 1997: An explicit flux-form semi-Lagrangian shallow-water model on the sphere. *Quart. J. Roy. Meteor. Soc.*, **123 (544)**, 2477–2498.

Lorenz, E. N., 1969: The predictability of a flow which possesses many scales of motion. *Tellus*, **21 (3)**, 289–307.

Lorenz, E. N., 1973: On the Existence of Extended Range Predictability. *J. Appl. Meteor.*, 543–546.

Mamalakis, A., I. Ebert-Uphoff, and E. A. Barnes, 2022: Neural network attribution methods for problems in geoscience: A novel synthetic benchmark dataset. *Environmental Data Science*, **1**, e8.

Mariotti, A., P. M. Ruti, and M. Rixen, 2018: Progress in subseasonal to seasonal prediction through a joint weather and climate community effort. *npj Climate Atmos. Sci.*, **1 (1)**, 4.

Mariotti, A., and Coauthors, 2020: Windows of opportunity for skillful forecasts subseasonal to seasonal and beyond. *Bull. Amer. Meteor. Soc.*, **101 (5)**, E608–E625.

Matheson, J. E., and R. L. Winkler, 1976: Scoring rules for continuous probability distributions. *Manage. Sci.*, **22 (10)**, 1087–1096.

Mayer, K. J., and E. A. Barnes, 2021: Subseasonal Forecasts of Opportunity Identified by an Explainable Neural Network. *Geophys. Res. Lett.*, **48 (10)**, e2020GL092 092.

Miles, E. R., C. M. Spillman, J. A. Church, and P. C. McIntosh, 2014: Seasonal prediction of global sea level anomalies using an ocean-atmosphere dynamical model. *Climate Dyn.*, **43**, 2131–2145.

Murphy, A. H., 1973: A new vector partition of the probability score. *J. Appl. Meteor. Climatol.*, **12 (4)**, 595–600.

Nakkiran, P., G. Kaplun, D. Kalimeris, T. Yang, B. L. Edelman, F. Zhang, and B. Barak, 2019: SGD on Neural Networks Learns Functions of Increasing Complexity. *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, Curran Associates Inc., Red Hook, NY, USA.

NASEM, 2016: *Next Generation Earth System Prediction: Strategies for Subseasonal to Seasonal Forecasts*. The National Academies Press, Washington, DC, https://doi.org/10.17226/21873.

Nix, D. A., and A. S. Weigend, 1994: Estimating the mean and variance of the target probability distribution. *Proceedings of 1994 IEEE international conference on neural networks (ICNN'94)*, IEEE, Vol. 1, 55–60.

O'Neill, B. C., and Coauthors, 2016: The Scenario Model Intercomparison Project (ScenarioMIP) for CMIP6. *Geosci. Model Dev.*, **9 (9)**, 3461–3482, https://doi.org/10.5194/gmd-9-3461-2016, URL https://gmd.copernicus.org/articles/9/3461/2016/.

Paszke, A., and Coauthors, 2019: Pytorch: an imperative style, high-performance deep learning library. *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, Curran Associates Inc., Red Hook, NY, USA.

Pedregosa, F., and Coauthors, 2011: Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.*, **12**, 2825–2830, URL https://jmlr.csail.mit.edu/papers/v12/pedregosa11a.html.

Penduff, T., M. Juza, L. Brodeau, G. C. Smith, B. Barnier, J.-M. Molines, A.-M. Treguier, and G. Madec, 2010: Impact of global ocean model resolution on sea-level variability with emphasis on interannual time scales. *Ocean Sci.*, **6 (1)**, 269–284.

Perezhogin, P., L. Zanna, and C. Fernandez-Granda, 2023: Generative Data-Driven Approaches for Stochastic Subgrid Parameterizations in an Idealized Ocean Model. *J. Adv. Model. Earth Syst.*, **15 (10)**, e2023MS003 681.

Piecuch, C., and R. Ponte, 2011: Mechanisms of interannual steric sea level variability. *Geophys. Res. Lett.*, **38 (15)**.

Qiu, B., 2002: Large-scale variability in the midlatitude subtropical and subpolar North Pacific Ocean: Observations and causes. *J. Phys. Oceanogr.*, **32 (1)**, 353–375.

Qu, Y., S. Jevrejeva, J. Williams, and J. C. Moore, 2022: Drivers for seasonal variability in sea level around the China seas. *Global Planet. Change*, **213**, 103 819.

Rao, S. A., and T. Yamagata, 2004: Abrupt termination of Indian Ocean dipole events in response to intraseasonal disturbances. *Geophys. Res. Lett.*, **31 (19)**.

Roberts, C., D. Calvert, N. Dunstone, L. Hermanson, M. Palmer, and D. Smith, 2016: On the Drivers and Predictability of Seasonal-to-Interannual Variations in Regional Sea Level. *J. Climate*, **29 (21)**, 7565–7585.

Rodgers, K., and Coauthors, 2021: Ubiquity of human-induced changes in climate variability. *Earth Syst. Dyn.*, **12 (4)**, 1393–1411.

Rossby, C.-G., 1939: Relation between variations in the intensity of the zonal circulation of the atmosphere and the displacements of the semi-permanent centers of action. *J. Mar. Res.*, **2**, 38–55.

Saji, N., B. N. Goswami, P. Vinayachandran, and T. Yamagata, 1999: A dipole mode in the tropical Indian Ocean. *Nature*, **401 (6751)**, 360–363.

Salmon, R., 1998: *Lectures on Geophysical Fluid Dynamics*. OUP USA.

Smith, R., and Coauthors, 2010: The Parallel Ocean Program (POP) Reference Manual. *LAUR-01853*, **141**, 1–140.

Srivastava, N., G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, 2014: Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *J. Mach. Learn. Res.*, **15 (1)**, 1929–1958.

Stuecker, M. F., A. Timmermann, F.-F. Jin, Y. Chikamoto, W. Zhang, A. T. Wittenberg, E. Widiasih, and S. Zhao, 2017: Revisiting ENSO/Indian Ocean dipole phase relationships. *Geophys. Res. Lett.*, **44 (5)**, 2481–2492.

Sundararajan, M., A. Taly, and Q. Yan, 2017: Axiomatic attribution for deep networks. *International Conference on Machine Learning*, PMLR, 3319–3328.

Vallis, G. K., 2017: *Atmospheric and Oceanic Fluid Dynamics*. Cambridge University Press.

Vitart, F., and Coauthors, 2017: The Subseasonal to Seasonal (S2S) Prediction Project Database. *Bull. Amer. Meteor. Soc.*, **98 (1)**, 163–173.

Wang, G., H.-L. Ren, J. Liu, and X. Long, 2023: Seasonal predictions of sea surface height in BCC-CSM1.1m and their modulation by tropical climate dominant modes. *Atmos. Res.*, **281**, 106 466.

Wang, H., R. Murtugudde, and A. Kumar, 2016: Evolution of Indian Ocean dipole and its forcing mechanisms in the absence of ENSO. *Climate Dynamics*, **47**, 2481–2500.

Webster, P. J., A. M. Moore, J. P. Loschnigg, and R. R. Leben, 1999: Coupled ocean–atmosphere dynamics in the Indian Ocean during 1997–98. *Nature*, **401 (6751)**, 356–360.

Wunsch, C., and D. Stammer, 1997: Atmospheric loading and the oceanic "inverted barometer" effect. *Rev. Geophys.*, **35 (1)**, 79–107.

Xie, S.-P., H. Annamalai, F. A. Schott, and J. P. McCreary, 2002: Structure and Mechanisms of South Indian Ocean Climate Variability. *J. Climate*, **15 (8)**, 864–878.

Yang, Y., S.-P. Xie, L. Wu, Y. Kosaka, N.-C. Lau, and G. A. Vecchi, 2015: Seasonality and predictability of the Indian Ocean dipole mode: ENSO forcing and internal variability. *Journal of Climate*, **28 (20)**, 8021–8036.