# Samudra: An AI Global Ocean Emulator for Climate

**Surya Dheeshjith[1], Adam Subel[1], Alistair Adcroft[2], Julius Busecke[4], Carlos Fernandez-Granda[1,3], Shubham Gupta[1], Laure Zanna[1]**

[1]Courant Institute of Mathematical Sciences, New York University
[2]Program in Atmospheric and Oceanic Sciences, Princeton University
[3]Center for Data Science, New York University
[4]Lamont Doherty Earth Observatory, Columbia University

**Key Points:**

- We develop a global, 3D, ocean autoregressive machine learning emulator for climate studies.
- The emulator, based on a UNet architecture, is stable for centuries, producing accurate climatologies and variability of ocean variables.
- The emulator training is robust to changes in seeds and initial conditions in the data.

Corresponding author: Surya Dheeshjith, `sd5313@nyu.edu`

**Abstract**

AI emulators for forecasting have emerged as powerful tools that can outperform conventional numerical predictions. The next frontier is to build emulators for long climate simulations with skill across a range of spatiotemporal scales, a particularly important goal for the ocean. Our work builds a skillful global emulator of the ocean component of a state-of-the-art climate model. We emulate key ocean variables, sea surface height, horizontal velocities, temperature, and salinity, across their full depth. We use a modified ConvNeXt UNet architecture trained on multi-depth levels of ocean data. We show that the ocean emulator – *Samudra* – which exhibits no drift relative to the truth, can reproduce the depth structure of ocean variables and their interannual variability. Samudra is stable for centuries and 150 times faster than the original ocean model. Samudra struggles to capture the correct magnitude of the forcing trends and simultaneously remains stable, requiring further work.

**Plain Language Summary**

AI tools are proving extremely effective in making fast and accurate predictions on weather to seasonal timescales. Capturing decadal to centennial changes, as those arising from ocean dynamics, remains an outstanding challenge for machine learning methods. We built an advanced AI model called "Samudra" to simulate global ocean behavior. Samudra is trained on simulated data from a state-of-the-art ocean climate model and predicts key ocean features such as sea surface height, currents, temperature, and salinity throughout the ocean's depth. Samudra can accurately recreate patterns in ocean variables, including year-to-year changes. It is stable over centuries and is 150 times faster than traditional ocean models. However, Samudra still faces challenges in balancing stability with accurately predicting the effects of external factors (like climate trends), and further improvements are needed to address this limitation.

## 1 Introduction

The recent success of emulators for components of the climate system, primarily the atmosphere, continues to produce remarkable outcomes, e.g., achieving state-of-the-art performance for weather prediction tasks (Kochkov et al., 2024; Bi et al., 2023; Price et al., 2023) and promising results reproducing climate models over decadal (Cachay et al., 2024) to multi-decadal timescales (Watt-Meyer et al., 2023).

Existing work on ocean emulation has mainly been limited to the surface and upper ocean, or to steady forcing. Several works focusing on surface ocean variables show results for timescales of years to a decade (Subel & Zanna, 2024; Dheeshjith et al., 2024; Gray et al., 2024). Emulators that include subsurface information have focused on the weekly to decadal timescales and at most the upper 1000m (Xiong et al., 2023; Guo et al., 2024; Holmberg et al., 2024), using a range of machine learning architectures (e.g., graph neural networks, Transformers). Bire et al. (2023) explored longer timescales within a simplified ocean model with idealized steady forcing. Finally, a first seasonal coupled atmosphere-ocean emulation has shown promising results Wang et al. (2024) considering the upper 300m of the ocean. These ocean and atmosphere emulators have focused on several tasks, from seasonal forecasts based on reanalysis data to building a surrogate of a numerical model for evaluation and prediction.

Building emulators (or surrogate models) of traditional numerical climate models aims to leverage the computational efficiency of machine learning approaches to reduce the often prohibitive computational cost of running a large number of numerical experiments with the original (usually CPU-based) climate model. One of the main benefits of building emulators is the ability to run large ensembles. For example, using large ensembles with different initial conditions, one can probe the likelihood of extreme events,

explore the climate response to a range of forcing scenarios (e.g., greenhouse gases), and enhance numerical model development by reducing the number of perturbed parameter experiments typically used for calibration (Maher et al., 2021; Mahesh et al., 2024). Emulators can be a useful tool for accelerating long spin-up integration or replacing full model components (Khatiwala, 2024). These emulators can also help with data assimilation, replacing an expensive numerical model with a fast surrogate to generate affordable ensembles or an approximate adjoint, to maintain accuracy with reduced cost (Manshausen et al., 2024).

Our goal here is to reproduce the full-depth ocean state for four 3D and one 2D prognostic variables, using a time-dependent realistic atmospheric forcing as input, extending the work of Subel and Zanna (2024); Dheeshjith et al. (2024). At rollout lengths of nearly a decade, our emulator shows considerable skill across several key diagnostics (mean and variance) when compared to the parent numerical model output, which is our ground truth. In particular, both the temperature structure as a function of depth and the El Niño-Southern Oscillation (ENSO) variability are well reproduced by the emulator.

Simultaneously capturing variables with vastly different timescales, such as velocity (that can contain fast fluctuations) and salinity (typically slow fluctuations), is an outstanding issue for long integrations (already encountered by Subel and Zanna (2024)). To alleviate this problem, we introduce an additional emulator by focusing on the thermodynamics variables (i.e. potential temperature and salinity only). This additional emulator captures the slowly varying changes in potential temperature and salinity on timescales of decades to centuries.

We show that our emulator can retain skill and remain stable for centuries for experiments equivalent to both control and climate change simulations. However, we also note that this stability is accompanied by a weak response to climate change forcing. This proof-of-concept work demonstrates (to our knowledge) the first ocean emulator capable of reproducing the full-depth (from the surface down to the ocean floor) ocean temperature structure and its variability, and running for multiple centuries in a realistic configuration with time-dependent forcing.

The paper is organized as follows. We discuss the data and all emulator details in Section 2. We explore the properties of the trained emulator across a test dataset and several multi-decadal experiments with a range of climate forcing in Section 3. We present our conclusions in Section 4.

## 2 Methods

We built an autoregressive ocean emulator from data generated by a state-of-the-art numerical ocean simulation. Below, we describe the data, the emulator, the architecture, and the training and evaluation of the emulator.

### 2.1 Data

The data was generated by OM4, (Adcroft et al., 2019), an ocean general circulation model used as the ocean component of the state-of-the-art coupled climate model CM4 (Held et al., 2019). The circulation model was initialized with hydrography from the World Ocean Atlas (Levitus et al., 2015) and forced with atmospheric reanalysis, following the OMIP-2 protocol, with version 1.4 of the JRA reanalysis (Tsujino et al., 2020). The model is run for 65 years (1958-2022).

The ocean prognostic variables are potential temperature ($\theta_O$), salinity ($S$), sea surface height (SSH), oceanic zonal ($u$), and meridional ($v$) velocity components. The circulation model has 75 degrees of freedom in the vertical for each 3D prognostic variable,

which we conservatively remap onto 19 fixed-depth levels of variable thickness - [2.5, 10, 22.5, 40, 65, 105, 165, 250, 375, 550, 775, 1050, 1400, 1850, 2400, 3100, 4000, 5000, 6000]m to reduce the data size. We also conservatively coarsen the data in time using a 5-day simple average in geopotential coordinates, averaging over the fastest waves resolved by the circulation model (which originally used a 20-minute time-step).

At this stage, the native horizontal grid for the data has a nominal resolution of $1/4°$ resolution but is curvilinear and has three poles (grid singularities) inland. We further post-process by filtering with an 18 by 18 cell Gaussian kernel using the gcm-filters package (Loose et al., 2022), and then conservatively interpolate onto a $1° \times 1°$ global geographic (latitude-longitude) grid using the xESMF package (Zhuang et al., 2023). Before the spatial conservative interpolation, we interpolate the velocities to the cell center using the xGCM package (Abernathey et al., 2022) and rotate the velocity vectors so that the $u$ and $v$ variables indicate purely zonal (east-west) and meridional (north-south) flow, respectively.

### 2.2 Ocean Emulator

The variables used to create the ocean emulator from the numerical model are as follows:

1. The ocean state $\mathbf{\Phi} = (\theta_O, S, \text{SSH}, u, v)$, which includes all 19 depth levels. We distinguish the set of thermodynamics variables as the subset consisting of $\mathbf{\Phi}_{\text{thermo}} = (\theta_O, S, \text{SSH})$, as opposed to the dynamic variables $\mathbf{\Phi}_{\text{dynamic}} = (u, v)$.
2. Atmosphere boundary conditions $\boldsymbol{\tau} = (\tau_u, \tau_v, Q, Q_{anom})$. This consists of the zonal surface ocean stress $\tau_u$, meridional surface ocean stress $\tau_v$, and net heat flux downward across the ocean surface Q (below the sea-ice) and its anomalies $Q_{anom}$. The net heat flux is a sum of the short- and long-wave radiative fluxes, sensible and latent heating, heat content of mass transfer, and heat flux due to frazil formation (see K4 and K5 of Griffies et al. (2016) for the precise definition of CMIP variable "hfds"). The heat flux anomalies are calculated by removing the climatological heat flux computed over the 65-year OM4 dataset.

Our emulator, $\mathcal{F}$, is built to autoregressively produce multiple future oceanic states given multiple previous oceanic states. Specifically, we use a 2-input - 2-output model configuration. Mathematically, we have,

$$\tilde{\mathbf{\Phi}}_{t+(n+1)\Delta t}, \tilde{\mathbf{\Phi}}_{t+(n+2)\Delta t} = \mathcal{F}(\tilde{\mathbf{\Phi}}_{t+(n-1)\Delta t}, \tilde{\mathbf{\Phi}}_{t+n\Delta t}, \boldsymbol{\tau}_{t+n\Delta t}) \tag{1}$$

where $n$ is a positive integer and $\tilde{\mathbf{\Phi}}$ represents the predicted ocean state by the emulator at time $t$. For the first time step, we use OM4 ocean states, $\mathbf{\Phi}_t$ and $\mathbf{\Phi}_{t-\Delta t}$, along with the corresponding atmospheric forcing, $\boldsymbol{\tau}_t$, to produce the first set of predictions. Subsequent ocean states are recursively produced by using the generated ocean states as input. We illustrate the rollout process of the emulator in Figure 1a). The use of multiple states provides additional context to the emulator, similar to the use of model time tendencies in PDE-based numerical integrations. In all of our experiments, $\Delta t = 5$ days.

### 2.3 Architecture

In this study, we rely on the ConvNeXt UNet architecture from (Dheeshjith et al., 2024) where the core blocks of a UNet (Ronneberger et al., 2015) are inspired by ConvNeXt blocks (Liu et al., 2022) adapted from (Karlbauer et al., 2023). The UNet implements downsampling based on average pooling and upsampling based on bilinear interpolation, which enables it to learn features at multiple scales. Each ConvNext block includes GeLU activations, increased dilation rates, and inverted channel bottlenecks. To save on computation, we did not use inverted channel depths and replaced the large $7 \times 7$ kernels with $3 \times 3$ kernels. We use batch normalization instead of layer normal-

ization, as it yielded better skill. The encoder and decoder consist of four ConvNeXt blocks, each with channel widths [200, 250, 300, 400]. The dilation rates used for both the encoder and decoder are [1, 2, 4, 8]. Additionally, we include a single ConvNext block (with channel width 400 and dilation 8) in the deepest section of the UNet before upsampling. The total number of parameters for the ConvNeXt UNet model used is 135M. We implemented periodic (or circular) padding in the longitudinal direction and zero padding at the poles as in (Dheeshjith et al., 2024).

The architecture is modified from Dheeshjith et al. (2024) to process the multiple depth level ocean data (as opposed to surface only). In the surface ocean emulator, which contains only a single depth level, each channel is associated with a variable. In the multi-depth ocean emulator, each channel is associated with a variable and a depth level. Our main emulator takes as input four 19-level oceanic variables ($\theta_O, S, u, v$), the surface variable SSH and four atmospheric boundary conditions ($\tau_u, \tau_v, Q, Q_{anom}$) and produces five output variables ($\theta_O, S, \mathrm{SSH}, u, v$). As discussed above, we use a 2-input 2-output model configuration and thus, there are $(4 \times 19 + 1) \times 2 + 4 = 158$ input and $(4 \times 19 + 1) \times 2 = 154$ output channels.

We train two emulators using the above architecture: (1) an emulator $\mathcal{F}_{\mathrm{thermo+dynamic}}$ that uses all the variables, $\mathbf{\Phi}_{\mathrm{thermo+dynamic}} = (\theta_O, S, \mathrm{SSH}, u, v)$, as input and output, and (2) an emulator $\mathcal{F}_{\mathrm{thermo}}$ that only uses the thermodynamic variables, $\mathbf{\Phi}_{\mathrm{thermo}} = (\theta_O, S, \mathrm{SSH})$.

### 2.4 Training Details

We illustrate the training of the model in Figure 1a). We train and validate the emulators using 2800 and 140 data samples corresponding to the years 1975 to 2012 and 2012 to 2014, respectively. Each sample is a 5-day mean of the full ocean state and atmospheric boundary conditions.

We ignore data over 1958-1975 due to the excessive model cooling while it adjusts from the warm initial conditions. This cooling does not reflect the forcing but rather an interior ocean model adjustment (see Sane et al. (2023) and S3). Note that some regions are still cooling post-1975 in this simulation, which biased some of our testing (see results).

The loss function used for optimization is

$$\mathcal{L}_t = \sum_{n=1}^{PN} \frac{1}{C\,Y\,X} \sum_{j=1}^{C} \sum_{k=1}^{Y} \sum_{l=1}^{X} \left( \tilde{\mathbf{\Phi}}_{t+n\Delta t}^{[j,k,l]} - \mathbf{\Phi}_{t+n\Delta t}^{[j,k,l]} \right)^2. \tag{2}$$

$\mathcal{L}_t$ is the total mean square error (MSE) loss function at time step t, where $P$ corresponds to the total number of input/output states used by the model in a single step, $N$ is the total number of recurrent passes, $C$, $Y$ and $X$ are the total number of output channels, height and width, respectively, of a single output state. Here, we set $P = 2$ to obtain a 2-input 2-output model configuration and $N = 4$ steps.

We use the Adam optimizer with a learning rate of $2e - 4$, which decays to zero using a Cosine scheduler. Our emulators are trained using 4 80GB A100 GPUs for 15 and 12 hours for the models $\mathcal{F}_{\mathrm{thermo+dynamic}}$ and $\mathcal{F}_{\mathrm{thermo}}$ respectively, with a total batch size of 16.

### 2.5 Evaluation

To evaluate the emulators, we take our initial conditions from 2014 and produce an 8-year rollout using the corresponding atmospheric forcing from 2014 to 2022. We compare the output from this rollout to held-out OM4 data to evaluate the emulator skill.
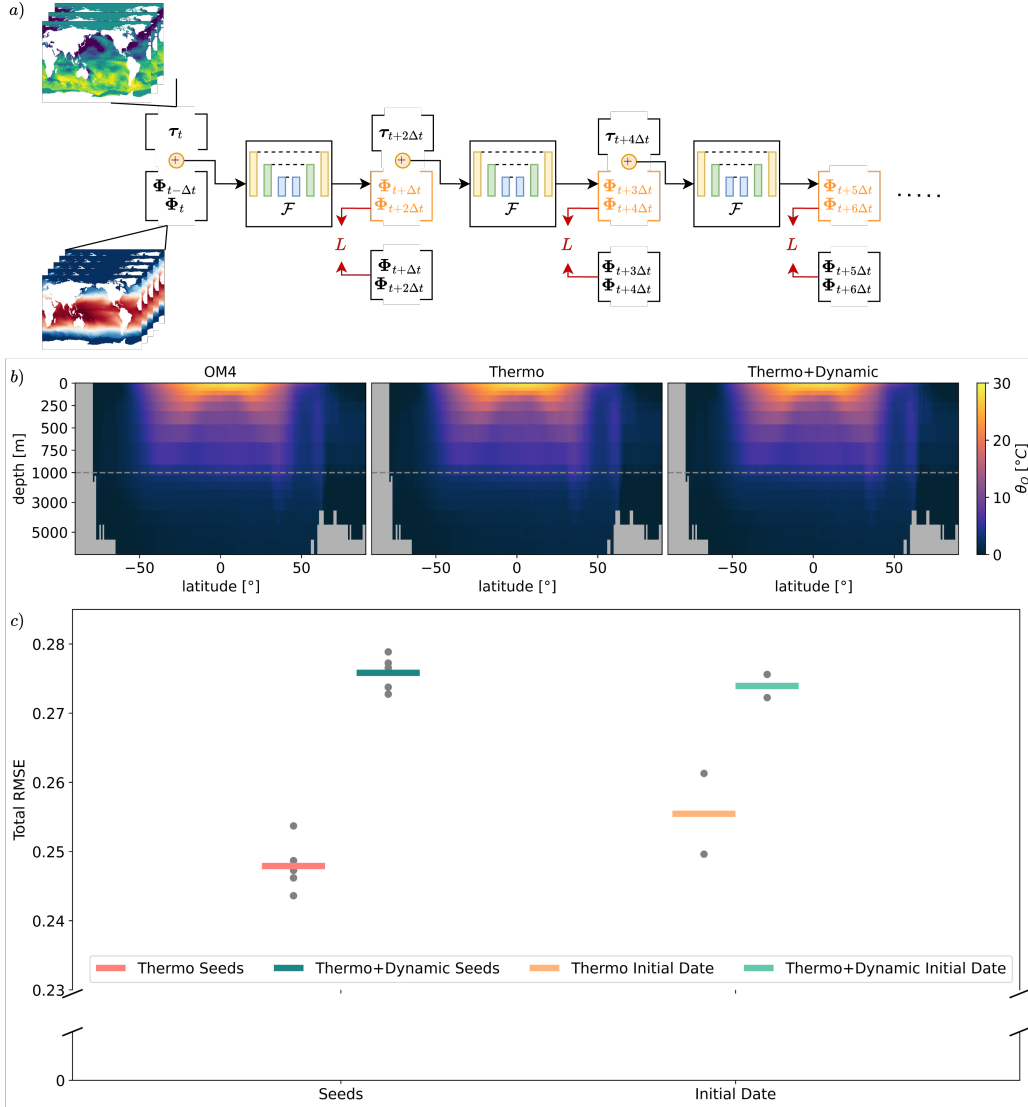
**Figure 1.** a) Schematic of the model training process, illustrating the mapping from input (ocean states and atmospheric forcing) to output (ocean states rolled out over several time steps). Initially, the ground truth ocean states, $\boldsymbol{\Phi}_t$ and $\boldsymbol{\Phi}_{t-\Delta t}$, along with the atmospheric forcing, $\boldsymbol{\tau}_t$, are provided as inputs to predict $\tilde{\boldsymbol{\Phi}}_{t+\Delta t}$ and $\tilde{\boldsymbol{\Phi}}_{t+2\Delta t}$. Predictions, along with ground truth atmospheric forcings, are then used as inputs for future steps in the unrolling process. b) Time-averaged potential temperature ($\theta_O$) depth-latitude profiles over the 8-year test set, comparing the ground truth OM4 (left) and predictions from $\mathcal{F}_{\text{thermo}}$ (middle) and $\mathcal{F}_{\text{thermo+dynamic}}$ (right). c) RMSE of 8-year test set predictions for different initial conditions of the emulators, $\mathcal{F}_{\text{thermo}}$ and $\mathcal{F}_{\text{thermo+dynamic}}$. Grey dots represent an RMSE instance of a single rollout, including runs from training on 5 unique model seeds per emulator and 2 additional rollouts initialized at states 6 months apart. Horizontal lines indicate the respective mean RMSE. RMSE is calculated over the common periods of each rollout.

In addition, we produce longer runs to assess the emulator's response, similar to control simulations, with arbitrarily long rollouts. The emulator is forced with atmospheric boundary conditions taken from 1990-2000, with a repeat 10-year cycle. This period is chosen specifically because it has a near-zero globally integrated heat flux forcing, which ensures minimal ocean drift. We also performed a 100-year and a 400-year control run (see SI).

For evaluations, we produce predictions using both $\mathcal{F}_{\text{thermo+dynamic}}$ and $\mathcal{F}_{\text{thermo}}$. All evaluations use a single 40GB A100 GPU. For each year of rollout, $\mathcal{F}_{\text{thermo+dynamic}}$ and $\mathcal{F}_{\text{thermo}}$ take about 90.52s and 47.2s, respectively. Thus, for the faster emulator, a century rollout takes approximately 1.3 hours. Roughly speaking, $\mathcal{F}_{\text{thermo}}$ takes about half the time to produce the same number of states in the rollout compared to $\mathcal{F}_{\text{thermo+dynamic}}$.

## 3 Results

### 3.1 Full-depth Global Ocean Emulator

We begin by evaluating the emulators $\mathcal{F}_{\text{thermo+dynamic}}$ and $\mathcal{F}_{\text{thermo}}$ against the ground truth to establish a baseline skill. Capturing the full-depth climatological profiles of potential temperature and salinity is a key target of ocean numerical climate models in general and, therefore, a key target for our ocean climate emulators. The structure of the zonal mean of potential temperature (Figure 1b) is captured by the two emulators, demonstrating significant skill at reproducing the profile from OM4 (see S6 for salinity structure). The average mean absolute error (MAE) is $5.7\times10^{-3}$ $^{\circ}C$ for $\mathcal{F}_{\text{thermo+dynamic}}$ and $4.5\times10^{-3}$ $^{\circ}C$ for $\mathcal{F}_{\text{thermo}}$, with a pattern correlation of roughly .99 for both emulators. The outputs show a robust thermocline structure, subtropical gyres, and a region of North Atlantic deep water formation.However, both emulators in the northern hemisphere show too warm and too salty high latitudes (around 55N), too cold and too fresh mid-latitudes, and Arctic signals down to 750m depth (Figures S2 and S7). The biases are consistent with underestimating the northward heat transport by the ocean, which is common to GCMs, including OM4. The potential temperature and salinity biases in the Southern Ocean for the $\mathcal{F}_{\text{thermo+dynamic}}$ emulator are reminiscent of residual transport changes, with opposite signed biases in the Southern Ocean and in the region north of it. The $\mathcal{F}_{\text{thermo}}$ emulator is warmer than the $\mathcal{F}_{\text{thermo+dynamic}}$, at most depths (Fig. S2).

We performed several experiments to test the sensitivity of the emulators to different training choices. The emulators' skill is unchanged when using different seeds and start dates, i.e., the trained models are statistically reproducible. We measure robustness by calculating the root mean square error (RMSE) of rollouts with 5 different seeds and rollouts initialized with ocean states taken 6 months apart. The RMSEs show little variance across the different trained models (Fig. 1c). The standard deviation of the RMSEs across training seeds in the emulators $\mathcal{F}_{\text{thermo}}$ and $\mathcal{F}_{\text{thermo+dynamic}}$ are 0.0033 and 0.00225, respectively.

The potential temperatures timeseries at 2.5m and 775m (Figure 2a) are further indicators that both emulators capture the climatological means and the upper ocean response to variable atmospheric forcing. The standard deviation of the 2.5m potential temperature for OM4, and the emulators $\mathcal{F}_{\text{thermo}}$ and $\mathcal{F}_{\text{thermo+dynamic}}$ are $6.8\times10^{-2}$ $^{\circ}C$, $4.35\times10^{-2}$ $^{\circ}C$ and $5.26\times10^{-2}$ $^{\circ}C$ respectively, while the standard deviations of the 775m potential temperature are $2.3\times10^{-3}$ $^{\circ}C$, $1.0\times10^{-3}$ $^{\circ}C$ and $2.1\times10^{-3}$ $^{\circ}C$, respectively. The standard deviations in each case are calculated after removing both the trend and the climatology from the timeseries (See Figure S8 for additional timeseries of potential temperature, along with salinity, zonal velocity, and meridional velocity, and Figure S10 for bias maps).

The emulator can skillfully emulate El Niño-Southern Oscillation (ENSO)' response in both warm and cold phases (Figure 2b) and S11). The smallest fluctuations in the

Nino 3.4 timeseries are the hardest for the emulators to capture. The emulator responses are in phase with OM4 for all years shown, but the amplitude is altered. $\mathcal{F}_{\text{thermo+dynamic}}$ exhibits higher skill than $\mathcal{F}_{\text{thermo}}$ in capturing the magnitude of ENSO events. We hypothesized that providing the velocities, whose data contain shorter time-scales and larger variability, helps the emulator produce larger ENSO events. $\mathcal{F}_{\text{thermo}}$ still manages to detect the correct phase and structure (Figure 2 b), d)) despite producing events with smaller magnitudes, both at the surface and in the upper ocean. The emulators capture the deepening and shoaling of the equatorial thermocline from equatorial Kelvin waves for the strongest events (Figure 2 d), e)). The magnitude of subsurface anomalies for the emulators is weaker than for OM4. Considering the Nino 3.4 timeseries (Figure 2 b), the MAE is 0.0077 $°C$ for $\mathcal{F}_{\text{thermo+dynamic}}$ and 0.0124 $°C$ for $\mathcal{F}_{\text{thermo}}$, with corresponding correlations of 0.905 and 0.7017, respectively. For the ENSO profiles (Figure 2 (c)-(e)), the MAE is 0.01 $°C$ and 0.07 $°C$ for the emulators $\mathcal{F}_{\text{thermo+dynamic}}$ and $\mathcal{F}_{\text{thermo}}$ respectively, and their corresponding pattern correlations are 0.976 and 0.973, respectively.

For the ocean emulator $\mathcal{F}_{\text{thermo+dynamic}}$ that uses all variables, we noticed that the potential temperature and salinity fields exhibit atypically high spatial variability, with scales more characteristic of velocity so we posit that this results from using velocity inputs. This result is consistent with Subel and Zanna (2024). We hypothesize that this may arise from the large separation in timescales and variability between velocity and potential temperature in the ocean.

Finally, despite capturing the mean and climatology of ocean variables, the emulators struggle to capture the magnitude of the small but systematic potential temperature trends (Figure S1 global mean $10^{-3}$ $°C/yr$) over the same 8-year period (Figure 2a and S1, S3); for most depths the trained models underestimate trends by 20% to 50% relative to OM4. Of the two emulators, $\mathcal{F}_{\text{thermo}}$ has higher skill in capturing the global heat changes (Figure S9). The salinity trends in OM4 are weak, due to the small forcing, and to the use of salinity restoring boundary conditions. For both emulators, the trends are 7-8 orders of magnitude less than the mean value, consistent with the numerical representation of variables within the learned models, and suggesting the models captured the conservation of properties inherent in the OM4 data without strict conservation being imposed (Figures S4-S5).

### 3.2 Long-term stability

We also evaluated, without retraining, the ability of the emulators to produce long control experiments. Specifically, for these experiments, we use repeat boundary conditions over 10 years (described in Section 2.5) chosen to contribute a near-zero net heat flux, allowing the emulators to run for arbitrarily long periods of time while minimizing potential temperature drift.

Both emulators converge to an equilibrium, maintaining a global mean potential temperature close to OM4 throughout a century of integration (Figure 3a). The global mean temperatures are 3.225 $°C/yr$ for $\mathcal{F}_{\text{thermo}}$ and 3.215 $°C/yr$ for $\mathcal{F}_{\text{thermo+dynamic}}$, compared to 3.219 $°C/yr$ for OM4. In addition, $\mathcal{F}_{\text{thermo+dynamic}}$ over-predicts the variability in potential temperature, likely extrapolating some fast dynamics via the velocities variables. This issue is exacerbated in the deeper layers of the ocean, which have little variability in the original dataset. The temperature structure is again well preserved for the long rollouts (Figure 3b), with different structures in potential temperature biases (S12) than for the 8-year test data (S2).

We examine the emulators' respective skill in reproducing variability over these long timescales. Since we are reusing the same 10-year cycle to drive the emulator, we expected some persistent features to appear when looking at a phenomenon such as the response to ENSO. Although both emulators can produce appropriate Nino 3.4 anomalies for the entire century rollout (Figure 3c) and S13), $\mathcal{F}_{\text{thermo+dynamic}}$ shows stronger peak-to-peak
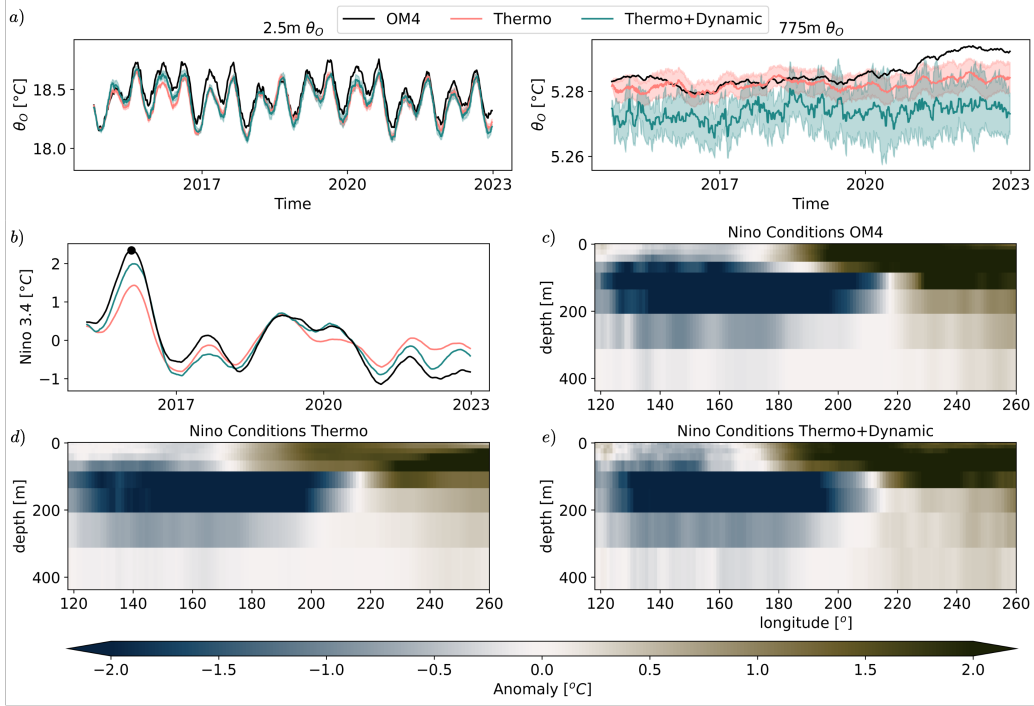
**Figure 2.** a) Spatially averaged timeseries of potential temperature $\theta_O$ at depths 2.5m (left) and 775m (right) over the test set comparing the ground truth OM4 (black), and predictions from $\mathcal{F}_{\text{thermo}}$ (red) and $\mathcal{F}_{\text{thermo+dynamic}}$ (green). The mean prediction and its variance (indicated by shading) are plotted over 5 initial seeds of training for each model. b) Nino 3.4 index timeseries over the test set for the ground truth (OM4, black) and predictions ($\mathcal{F}_{\text{thermo}}$, red; $\mathcal{F}_{\text{thermo+dynamic}}$, green). Anomalies are averaged over rolling 150-day windows. c-e) Meridionally averaged depth profile of potential temperature anomalies in the tropics during the peak Nino event (marked by a black dot in the timeseries) over the test set for OM4 (c), $\mathcal{F}_{\text{thermo}}$ (d) and $\mathcal{F}_{\text{thermo+dynamic}}$ (e). Anomalies in (c)-(e) are averaged over a 15-day window.

amplitude, but little cycle-to-cycle variability - perhaps due to the strong coupling of velocity with the wind stress forcing, whereas $\mathcal{F}_{\text{thermo}}$ shows more aperiodic variability across years.

To further test stability, we generate a 400-year rollout, with an identical forcing setup as for the century-long run. Both emulators remain stable (Figure S15). $\mathcal{F}_{\text{thermo}}$ has the added benefit of exhibiting long-term aperiodic variability in potential temperature and salinity, despite the repeat forcing, across the centuries.

## 4 Discussion

We produce a computationally cheap machine-learning (ML) emulator of a state-of-the-art ocean model, namely OM4 (Adcroft et al., 2019). The ML architecture consists of a modified ConvNeXt UNet (Dheeshjith et al., 2024). The reduced order model – *Samudra* – predicts key ocean variables, sea surface height, temperature, and salinity, across the full depth of the world oceans while remaining stable for centuries. Integrating OM4 for 100 years takes approximately 8 days using 4,671 CPU cores, whereas our fastest (thermo) emulator completes the same task in about 1.3 hours on a single 40GB A100 GPU. This represents approximately a 150x increase in SYPD (simulated years per day) for Samudra compared to OM4.

The emulator performs well on a range of metrics related to the model climatology and its variability on the test set and long control simulations. The emulator produces accurate climatologies over the last 8 years of the OM4 simulations and is robust to changes in seeds and initial conditions. Furthermore, it can capture variability (e.g., ENSO response to forcing). Therefore, these emulators could be used to study the contemporary ocean and climate at a significant reduction in cost compared to OM4.

The emulators, however, struggle to capture trends under a range of surface heat flux forcing (see Supporting Information), similarly to Dheeshjith et al. (2024) for surface emulators. We performed idealized forced experiments using the same repeated atmospheric forcing generated for the control experiment and a spatially uniform linear forcing of varying magnitudes for the surface heat flux. Figure S16 showcases the ocean heat content trends predicted by $\mathcal{F}_{\text{thermo}}$ under linear surface heat flux increases of 1, 0.5, 0.25, and 0 $W/m^2$. The patterns of ocean heat uptake are reminiscent of ocean-only and coupled forced numerical experiments (Todd et al., 2020; Couldrey et al., 2020), with dipole patterns in the Southern Ocean and North Atlantic sinking region (Figure S14). However, the magnitude of change is too weak compared to the forcing (Figure S16). Similar behavior of weak generalization under climate change is also observed in the atmosphere climate emulator, ACE (Watt-Meyer et al., 2023) but improved when a slab ocean model is added (Clark et al., 2024).

In our work, we could not produce an emulator that would capture the trends in the test data and remain stable for multiple centuries. Further work is needed to thoroughly explore the reasons for the issues and would require new numerical simulations, which are beyond the scope of this work. However, we outline some hypotheses. Specifically, there are several possible reasons for the lack of generalization associated with the weak warming trends: the ocean data for training, the atmospheric forcing data, the model formalism, or the machine learning architecture.

The ocean training data is from the OM4 run, similar to Sane et al. (2023). The effects of an initial drift can be alleviated by pruning years 1958 to 1975 from the training data, which removes the bulk of this adjustment period. Yet, different depths and regions adjust more slowly, and some of this continued adjustment may remain in the data since the time scale of equilibration of the model is 100's of years. The other reason for the trend bias could possibly come from the forcing datasets. The atmospheric forcing imposed on the ocean implicitly results from the real ocean-atmosphere coupling.
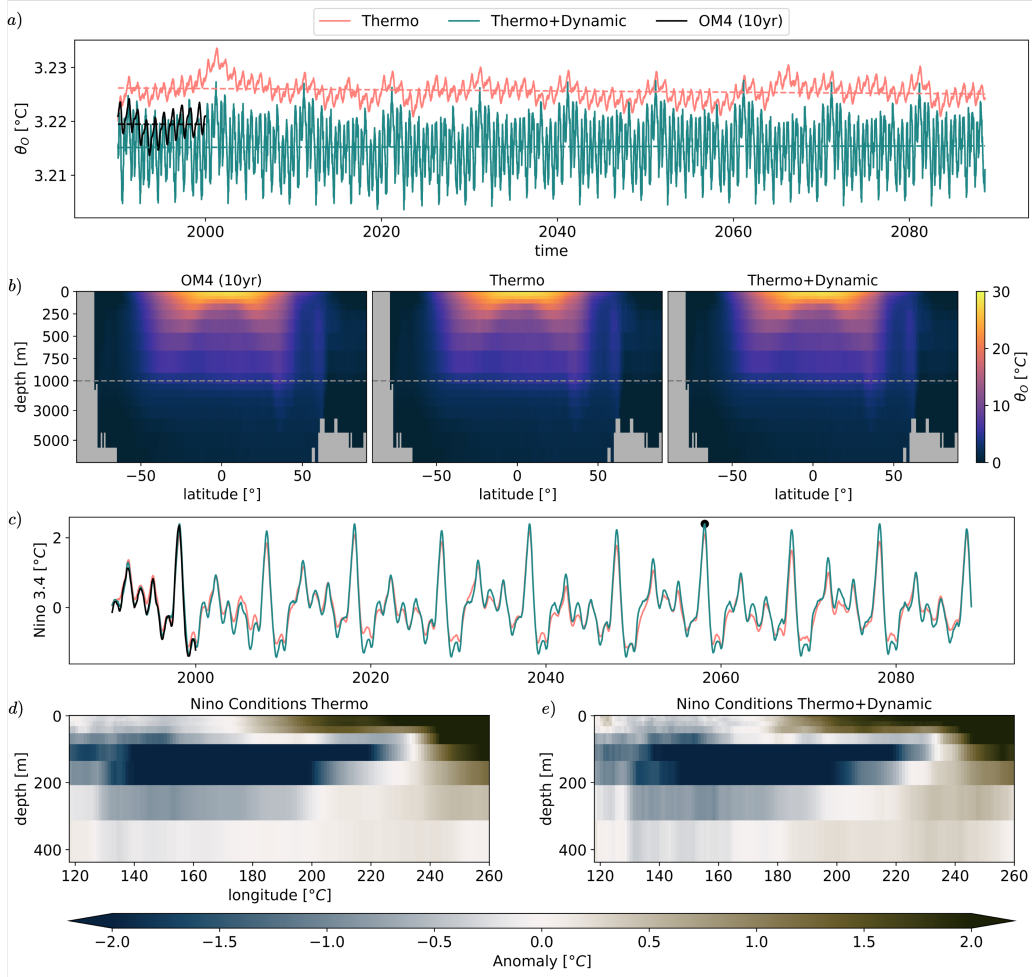
**Figure 3.** a) Globally averaged potential temperature ($\theta_O$) timeseries over a 100-year control run, comparing the 10-year ground truth OM4 (black) and predictions from $\mathcal{F}_{\text{thermo}}$ (red) and $\mathcal{F}_{\text{thermo+dynamic}}$ (green). b) Time-averaged potential temperature ($\theta_O$) depth profile over a 100-year control run, comparing the 10-year ground truth OM4 (left) and predictions from $\mathcal{F}_{\text{thermo}}$ (middle) and $\mathcal{F}_{\text{thermo+dynamic}}$ (right). c) Nino 3.4 index timeseries over a 100-year control run, comparing the 10-year repeat for the ground truth (OM4, black) and predictions ($\mathcal{F}_{\text{thermo}}$, red; $\mathcal{F}_{\text{thermo+dynamic}}$, green). d-e) Meridionally averaged depth profile of potential temperature anomalies in the tropics during the peak Nino event (marked by a black dot in the timeseries) over the test set for $\mathcal{F}_{\text{thermo}}$ (d) and $\mathcal{F}_{\text{thermo+dynamic}}$ (e). Anomalies are as in Fig. 2.

Therefore, the atmospheric forcing has felt a changing ocean circulation, particularly in the North Atlantic (Chemke et al., 2020). The resulting effect is that the "forcing" applied to the ocean emulator is not entirely decoupled from the ocean response, potentially leading to some biases in the response as in Todd et al. (2020); Couldrey et al. (2020); Zanna et al. (2019). We alleviated these issues by adding an extra forcing input, namely the cumulative heat forcing, which led to a more skillful model capable of capturing the global warming trend. However, this model was unstable under climate change forcing past 50 years. Another possibility is learning to predict the model state directly might not be optimal. We explored learning tendencies which improved performance for the warming trends but, again, was unstable over long timescales. Finally, different architectures might also play a role. However, our original exploration (Dheeshjith et al., 2024) has not yet shown much improvement over the ConvNext architecture. Therefore, we focused on the stable emulator in this paper as described above. A significant challenge going forward is designing faithful emulators capable of capturing trends while remaining stable in long rollouts.

Despite the inability to respond to future climate forcing, Samudra is skillful at emulating the contemporary ocean. Without further modification, Samudra could be used, as is, in studies requiring large ensembles (e.g., uncertainty quantification, extreme events) or to enhance and accelerate operational applications (e.g., data assimilation). More opportunities emerge if we consider refining training for Samudra, e.g., to revised versions of OM4 or to other models, which could greatly accelerate climate model development by allowing evaluations of long, yet affordable, rollouts. Samudra is more than a proof of concept for affordable emulation of expensive ocean circulation models, and could be used off-the-shelf for many applications.

## Open Research Section

The code for generating rollouts and plots is available on GitHub at https://github.com/m2lines/Samudra, while the model weights and data are hosted on Hugging Face at https://huggingface.co/M2LInES/Samudra and https://huggingface.co/datasets/M2LInES/Samudra_OM4, respectively.

The software used for training the models will be placed on GitHub at the revision stage of this work. For publication the code will be version tagged and archived via zenodo.

## References

Abernathey, R. P., Busecke, J. J. M., Smith, T. A., Deauna, J. D., Banihirwe, A., Nicholas, T., . . . Thielen, J. (2022, November). *xgcm.* Zenodo. Retrieved from `https://doi.org/10.5281/zenodo.7348619` doi: 10.5281/zenodo.7348619

Adcroft, A., Anderson, W., Balaji, V., Blanton, C., Bushuk, M., Dufour, C. O., . . . Zhang, R. (2019). The GFDL Global Ocean and Sea Ice Model OM4.0: Model Description and Simulation Features. *Journal of Advances in Modeling Earth Systems*, *11*(10), 3167–3211. Retrieved from `https://onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001726` doi:

10.1029/2019MS001726

Bi, K., Xie, L., Zhang, H., Chen, X., Gu, X., & Tian, Q. (2023). Accurate medium-range global weather forecasting with 3d neural networks. *Nature*, *619*(7970), 533–538.

Bire, S., Lütjens, B., Azizzadenesheli, K., Anandkumar, A., & Hill, C. N. (2023). Ocean emulation with fourier neural operators: Double gyre. *Authorea Preprints*.

Cachay, S. R., Henn, B., Watt-Meyer, O., Bretherton, C. S., & Yu, R. (2024). Probabilistic emulation of a global climate model with spherical dyffusion. *arXiv preprint arXiv:2406.14798*.

Chemke, R., Zanna, L., & Polvani, L. M. (2020). Identifying a human signal in the north atlantic warming hole. *Nature communications*, *11*(1), 1540.

Clark, S. K., Watt-Meyer, O., Kwa, A., McGibbon, J., Henn, B., Perkins, W. A., . . . Harris, L. M. (2024). Ace2-som: Coupling to a slab ocean and learning the sensitivity of climate to changes in co 2. *arXiv preprint arXiv:2412.04418*.

Couldrey, M. P., Gregory, J. M., Dias, F. B., Dobrohotoff, P., Domingues, C. M., Garuba, O., . . . others (2020). What causes the spread of model projections of ocean dynamic sea-level change in response to greenhouse gas forcing? *Climate Dynamics*, 1–33.

Dheeshjith, S., Subel, A., Gupta, S., Adcroft, A., Fernandez-Granda, C., Busecke, J., & Zanna, L. (2024). Transfer learning for emulating ocean climate variability across co 2 forcing. *arXiv preprint arXiv:2405.18585*.

Gray, M. A., Chattopadhyay, A., Wu, T., Lowe, A., & He, R. (2024). Long-term prediction of the gulf stream meander using oceannet: a principled neural operator-based digital twin. *EGUsphere*, *2024*, 1–23.

Griffies, S. M., Danabasoglu, G., Durack, P. J., Adcroft, A. J., Balaji, V., Böning, C. W., . . . Yeager, S. G. (2016, September). OMIP contribution to CMIP6: experimental and diagnostic protocol for the physical component of the Ocean Model Intercomparison Project. *Geoscientific Model Development*, *9*(9), 3231–3296. doi: https://doi.org/10.5194/gmd-9-3231-2016

Guo, Z., Lyu, P., Ling, F., Luo, J.-J., Boers, N., Ouyang, W., & Bai, L. (2024). Orca: A global ocean emulator for multi-year to decadal predictions. *arXiv preprint arXiv:2405.15412*.

Held, I. M., Guo, H., Adcroft, A., Dunne, J. P., Horowitz, L. W., Krasting, J., . . . Zadeh, N. (2019). Structure and Performance of GFDL's CM4.0 Climate Model. *Journal of Advances in Modeling Earth Systems*, *11*(11), 3691–3727. Retrieved from `https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001829` doi: 10.1029/2019MS001829

Holmberg, D., Clementi, E., & Roos, T. (2024). Regional ocean forecasting with hierarchical graph neural networks. *arXiv preprint arXiv:2410.11807*.

Karlbauer, M., Cresswell-Clay, N., Durran, D. R., Moreno, R. A., Kurth, T., & Butz, M. V. (2023). Advancing parsimonious deep learning weather prediction using the healpix mes. *Authorea Preprints*.

Khatiwala, S. (2024). Efficient spin-up of earth system models using sequence acceleration. *Science Advances*, *10*(18), eadn2839.

Kochkov, D., Yuval, J., Langmore, I., Norgaard, P., Smith, J., Mooers, G., . . . others (2024). Neural general circulation models for weather and climate. *Nature*, 1–7.

Levitus, S., Boyer, T., Garcia, H., Locarnini, R., Zweng, M., Mishonov, A., . . . Seidov, D. (2015). *World ocean atlas 2013 (NCEI accession 0114815).* doi: 10.7289/v5f769gt

Liu, Z., Mao, H., Wu, C.-Y., Feichtenhofer, C., Darrell, T., & Xie, S. (2022). A convnet for the 2020s. In *Proceedings of the ieee/cvf conference on computer vision and pattern recognition* (pp. 11976–11986).

Loose, N., Abernathey, R., Grooms, I., Busecke, J., Guillaumin, A., Yankovsky,

E., . . . Martin, P. (2022). Gcm-filters: A python package for diffusion-based spatial filtering of gridded data. *Journal of Open Source Software*, *7*(70), 3947. Retrieved from `https://doi.org/10.21105/joss.03947` doi: 10.21105/joss.03947

Maher, N., Milinski, S., & Ludwig, R. (2021). Large ensemble climate model simulations: introduction, overview, and future prospects for utilising multiple types of large ensemble. *Earth System Dynamics*, *12*(2), 401–418.

Mahesh, A., Collins, W., Bonev, B., Brenowitz, N., Cohen, Y., Elms, J., . . . others (2024). Huge ensembles part i: Design of ensemble weather forecasts using spherical fourier neural operators. *arXiv preprint arXiv:2408.03100*.

Manshausen, P., Cohen, Y., Pathak, J., Pritchard, M., Garg, P., Mardani, M., . . . Brenowitz, N. (2024). Generative data assimilation of sparse weather station observations at kilometer scales. *arXiv preprint arXiv:2406.16947*.

Price, I., Sanchez-Gonzalez, A., Alet, F., Andersson, T. R., El-Kadi, A., Masters, D., . . . others (2023). Gencast: Diffusion-based ensemble forecasting for medium-range weather. *arXiv preprint arXiv:2312.15796*.

Ronneberger, O., Fischer, P., & Brox, T. (2015). U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–miccai 2015: 18th international conference, munich, germany, october 5-9, 2015, proceedings, part iii 18* (pp. 234–241).

Sane, A., Reichl, B. G., Adcroft, A., & Zanna, L. (2023). Parameterizing vertical mixing coefficients in the ocean surface boundary layer using neural networks. *Journal of Advances in Modeling Earth Systems*, *15*(10), e2023MS003890.

Subel, A., & Zanna, L. (2024). Building ocean climate emulators. *arXiv preprint arXiv:2402.04342*.

Todd, A., Zanna, L., Couldrey, M., Gregory, J., Wu, Q., Church, J. A., . . . others (2020). Ocean-only fafmip: Understanding regional patterns of ocean heat content and dynamic sea level change. *Journal of Advances in Modeling Earth Systems*, *12*(8), e2019MS002027.

Tsujino, H., Urakawa, L. S., Griffies, S. M., Danabasoglu, G., Adcroft, A. J., Amaral, A. E., . . . Yu, Z. (2020, August). Evaluation of global ocean–sea-ice model simulations based on the experimental protocols of the Ocean Model Intercomparison Project phase 2 (OMIP-2). *Geoscientific Model Development*, *13*(8), 3643–3708. Retrieved from `https://gmd.copernicus.org/articles/13/3643/2020/` doi: https://doi.org/10.5194/gmd-13-3643-2020

Wang, C., Pritchard, M. S., Brenowitz, N., Cohen, Y., Bonev, B., Kurth, T., . . . Pathak, J. (2024). Coupled ocean-atmosphere dynamics in a machine learning earth system model. *arXiv preprint arXiv:2406.08632*.

Watt-Meyer, O., Dresdner, G., McGibbon, J., Clark, S. K., Henn, B., Duncan, J., . . . others (2023). Ace: A fast, skillful learned global atmospheric model for climate prediction. *arXiv preprint arXiv:2310.02074*.

Xiong, W., Xiang, Y., Wu, H., Zhou, S., Sun, Y., Ma, M., & Huang, X. (2023). Ai-goms: Large ai-driven global ocean modeling system. *arXiv preprint arXiv:2308.03152*.

Zanna, L., Khatiwala, S., Gregory, J. M., Ison, J., & Heimbach, P. (2019). Global reconstruction of historical ocean heat storage and transport. *Proceedings of the National Academy of Sciences*, *116*(4), 1126–1131.

Zhuang, J., raphael dussin, Huard, D., Bourgault, P., Banihirwe, A., Raynaud, S., . . . Li, X. (2023, September). *pangeo-data/xesmf: v0.8.2*. Zenodo. Retrieved from `https://doi.org/10.5281/zenodo.8356796` doi: 10.5281/zenodo.8356796

**Supporting Information**

**Text S1.** Here we describe how we calculated $Q_{anom}$.

$$Q_{anom}(t, y, x) = Q(t, y, x) - Clim(Q)(t, y, x) \tag{3}$$

where Clim is the climatology of Q over the entire data.

**Text S2.** Calculation of Metrics

Consider a predicted ocean state $\tilde{\mathbf{\Phi}}_t^{[j,k,l]}$, its corresponding ground truth state $\mathbf{\Phi}_t^{[j,k,l]}$ at time $t$, channel $j$, latitude $k$ and longitude $l$, and the normalized volume $V(j, k, l)$ at channel $j$, latitude $k$ and longitude $l$.

$$RMSE(\tilde{\mathbf{\Phi}}, \mathbf{\Phi}) = \frac{1}{T} \sum_t \sqrt{\sum_{j,k,l} V(j, k, l) \left( \tilde{\mathbf{\Phi}}_t^{[j,k,l]} - \mathbf{\Phi}_t^{[j,k,l]} \right)^2} \tag{4}$$

$$MAE(\tilde{\mathbf{\Phi}}, \mathbf{\Phi}) = \frac{1}{T} \sum_t \left| \sum_{j,k,l} V(j, k, l) \left( \tilde{\mathbf{\Phi}}_t^{[j,k,l]} - \mathbf{\Phi}_t^{[j,k,l]} \right) \right| \tag{5}$$

$$Corr(\tilde{\mathbf{\Phi}}, \mathbf{\Phi}) = \frac{1}{T} \sum_t \frac{\sum_{j,k,l} V(j, k, l) \tilde{\mathbf{\Phi}}_t^{[j,k,l]} \mathbf{\Phi}_t^{[j,k,l]}}{\sqrt{\sum_{j,k,l} V(j, k, l) (\tilde{\mathbf{\Phi}}_t^{[j,k,l]})^2 \sum_{j,k,l} V(j, k, l) (\mathbf{\Phi}_t^{[j,k,l]})^2}} \tag{6}$$

where $T$ is the time period over which we calculate the metrics.

**Figure S1.** Spatially averaged potential temperature ($\theta_O$) time series over an 8-year test set comparing the ground truth OM4 (black), and predictions from $\mathcal{F}_{\text{thermo}}$ (red), and $\mathcal{F}_{\text{thermo+dynamic}}$ (green). The mean prediction and its variance (indicated by shading) are plotted over 5 initial seeds of training for each model.
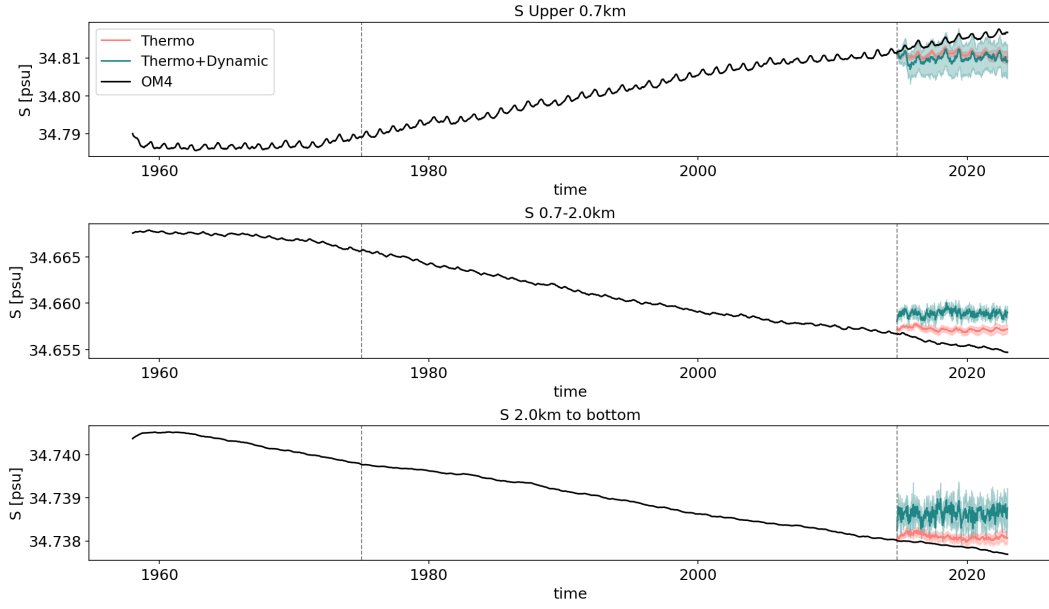


**Figure S2.** Time- and zonally-averaged potential temperature ($\theta_O$) biases (relative to OM4) for an 8-year test set: $\mathcal{F}_{\text{thermo}}$ (left), $\mathcal{F}_{\text{thermo+dynamic}}$ (center), and the difference between $\mathcal{F}_{\text{thermo}}$ and $\mathcal{F}_{\text{thermo+dynamic}}$ (right).



**Figure S3.** Spatially averaged potential temperature ($\theta_O$) trends of the entire ground truth data OM4 (black), and 8-year test set predictions from $\mathcal{F}_{\text{thermo}}$ (red) and $\mathcal{F}_{\text{thermo+dynamic}}$ (green) at depth levels 0–700 m, 700–2000 m, and 2000–6000 m. Vertical lines indicate the section of training data considered. The mean prediction and its variance (indicated by shading) are plotted over 5 initial seeds of training for each model.
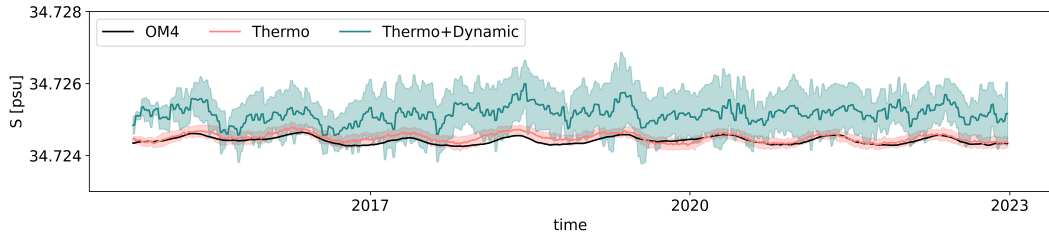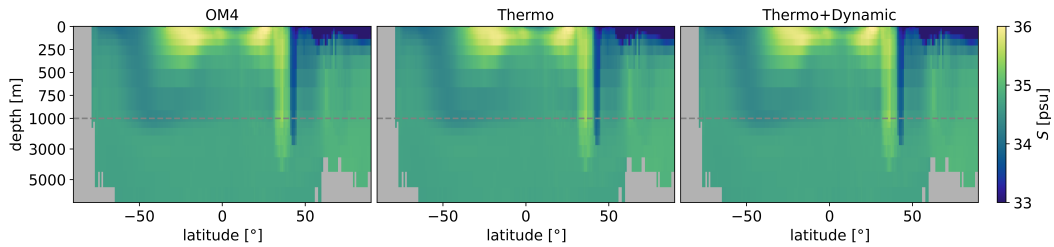
**Figure S4.** Spatially averaged salinity ($S$) trends of the entire ground data truth OM4 (black), and 8-year test set predictions from $\mathcal{F}_{\text{thermo}}$ (red) and $\mathcal{F}_{\text{thermo+dynamic}}$ (green) at depth levels 0–700 m, 700–2000 m, and 2000–6000 m. Vertical lines indicate the section of training data considered. The mean prediction and its variance (indicated by shading) are plotted over 5 initial seeds of training for each model.



**Figure S5.** Spatially averaged salinity ($S$) time series over an 8-year test set comparing the ground truth OM4 (black), and predictions from $\mathcal{F}_{\text{thermo}}$ (red), and $\mathcal{F}_{\text{thermo+dynamic}}$ (green). The mean prediction and its variance (indicated by shading) are plotted over 5 initial seeds of training for each model.



**Figure S6.** Time- and zonally-averaged salinity ($S$) for an 8-year test set: ground truth OM4 (left), $\mathcal{F}_{\text{thermo}}$ (center), and $\mathcal{F}_{\text{thermo+dynamic}}$ (right).

**Figure S7.** Time- and zonally-averaged salinity ($S$) biases (relative to OM4) for $\mathcal{F}_{\text{thermo}}$ (left), $\mathcal{F}_{\text{thermo+dynamic}}$ (center), and the difference between $\mathcal{F}_{\text{thermo}}$ and $\mathcal{F}_{\text{thermo+dynamic}}$ (right) for an 8-year test set.
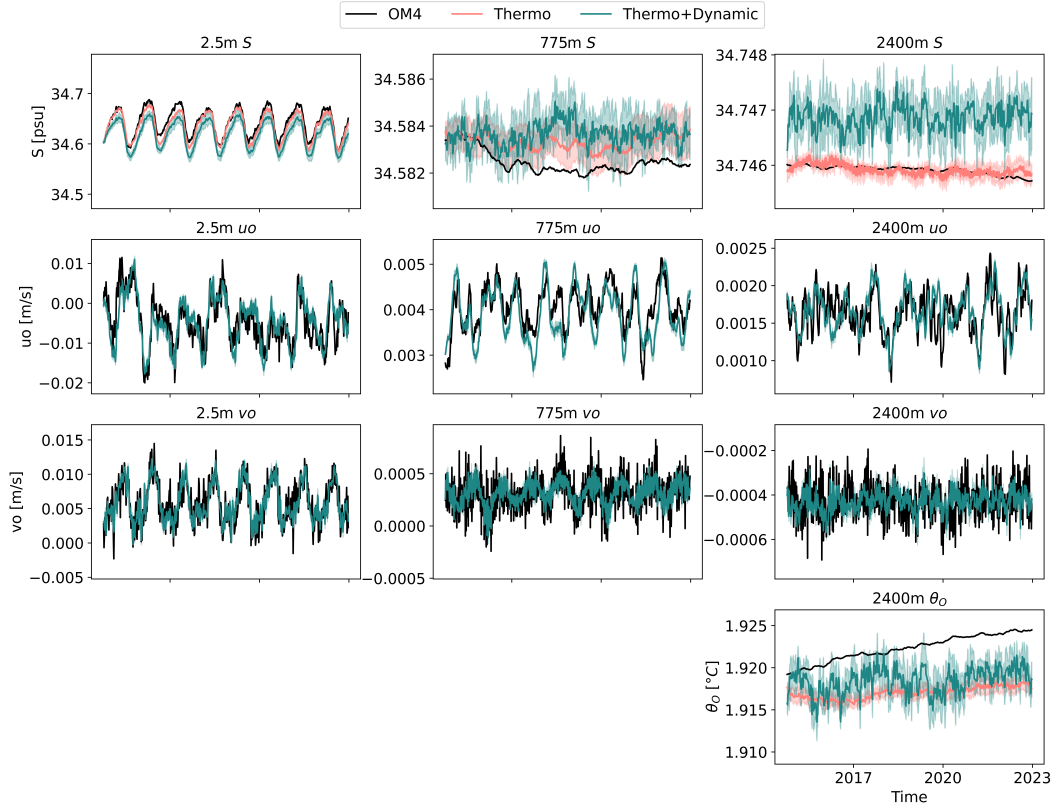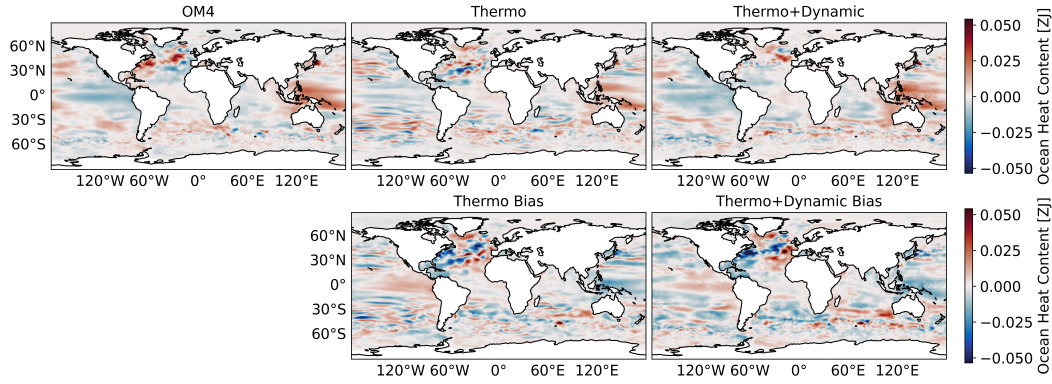


**Figure S8.** Spatially averaged time series over an 8-year test set for the ground truth OM4 (black), $\mathcal{F}_{\text{thermo}}$ (red), and $\mathcal{F}_{\text{thermo+dynamic}}$ (green). The first, second, and third rows correspond to salinity ($S$), zonal velocity ($uo$), and meridional velocity ($vo$) at depths of 2.5m, 775m, and 2400m, respectively. The final plot in the bottom row represents potential temperature ($\theta_O$) at 2400m. The mean prediction and its variance (indicated by shading) are plotted over 5 initial seeds of training for each model.

**Figure S9.** Global maps of Ocean Heat Content (OHC) evaluated over an 8-year test set, displaying the difference between the last and first year for the ground truth OM4 (top left), $\mathcal{F}_{\text{thermo}}$ (top center), and $\mathcal{F}_{\text{thermo+dynamic}}$ (top right). The corresponding bias maps are shown in the bottom row.
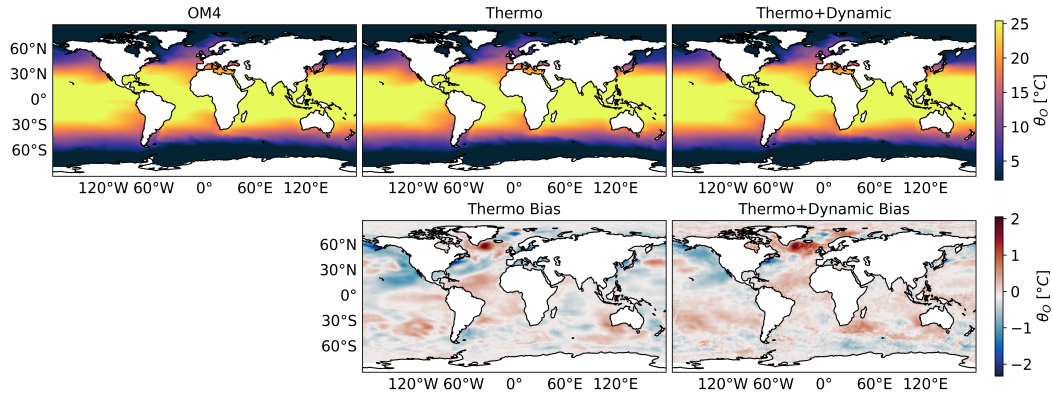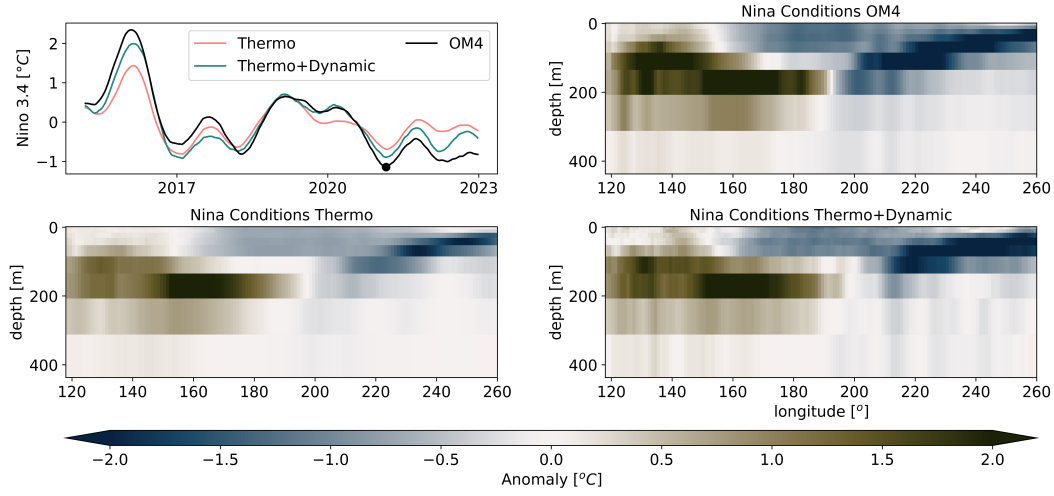


**Figure S10.** Time-averaged global maps of 2.5m potential temperature ($\theta_O$) evaluated over an 8-year test set for the ground truth OM4 (top left), $\mathcal{F}_{\text{thermo}}$ (top center), and $\mathcal{F}_{\text{thermo+dynamic}}$ (top right), with corresponding bias maps displayed in the bottom row.

**Figure S11.** Timeseries of Nino 3.4 index over an 8-year test set, comparing the ground truth OM4 (black) with predictions from $\mathcal{F}_{\text{thermo}}$ (red) and $\mathcal{F}_{\text{thermo+dynamic}}$ (green). Here, we consider the 2.5m temperature anomalies. Anomalies are calculated relative to the 8-year climatology of OM4 and each emulator. Additionally, the depth structure of anomalies is shown for the peak Nina event (marked by a black dot in the timeseries). Anomalies are averaged over rolling 150-day windows in the timeseries while the anomalies in the depth structures are averaged over a 15-day (3-snapshot) window to reduce mesoscale variability.



**Figure S12.** Time- and zonally-averaged potential temperature ($\theta_O$) biases (relative to OM4) for a 100-year control run forced with repeated atmospheric conditions taken from 1990-2000: $\mathcal{F}_{\text{thermo}}$ (left), $\mathcal{F}_{\text{thermo+dynamic}}$ (center), and the difference between $\mathcal{F}_{\text{thermo}}$ and $\mathcal{F}_{\text{thermo+dynamic}}$ (right). We compare the average for the 10-year period (1990–2000) of OM4 with the average of the 100-year emulator run.
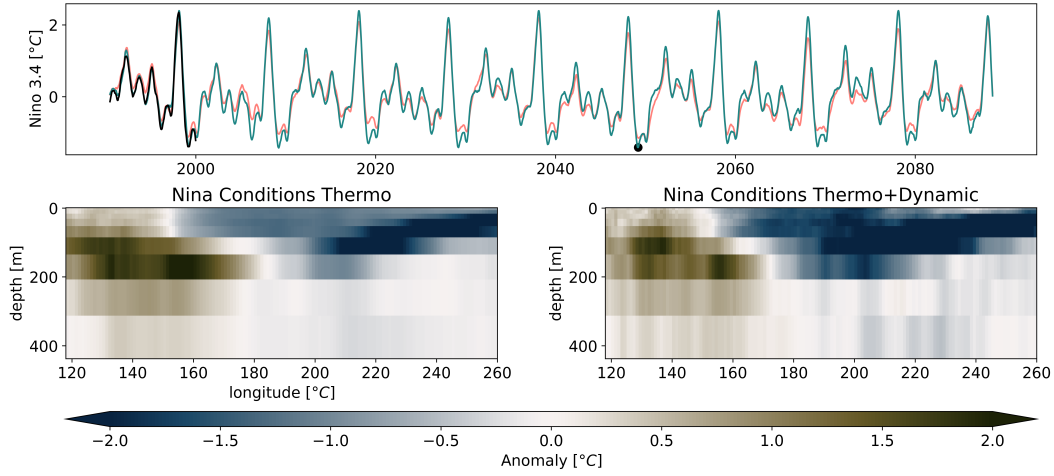
**Figure S13.** Timeseries of Nino 3.4 index over a 100-year control run, comparing the 10-year repeat ground truth OM4 (black) with predictions from $\mathcal{F}_{\text{thermo}}$ (red) and $\mathcal{F}_{\text{thermo+dynamic}}$ (green). Here, we consider the 2.5m temperature anomalies. Anomalies are calculated relative to the 10-year climatology of OM4 and 100-year climatology of each emulator. Additionally, the depth structure of anomalies is shown for the peak Nina event (marked by a black dot in the timeseries). Anomalies are averaged over rolling 150-day windows in the timeseries while the anomalies in the depth structures are averaged over a 15-day (3-snapshot) window to reduce mesoscale variability.
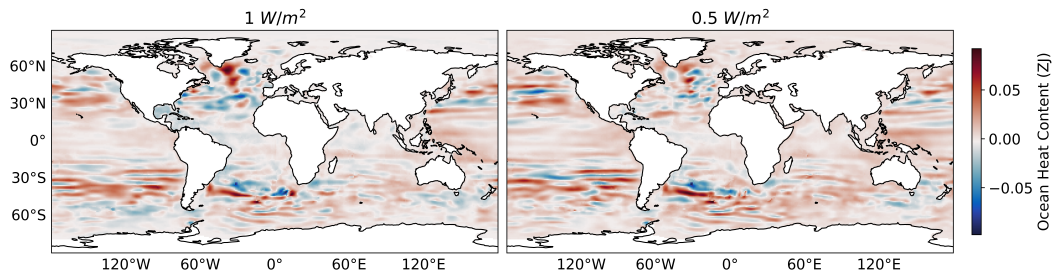


**Figure S14.** OHC Global Maps for the $\mathcal{F}_{\text{thermo}}$ emulator, evaluated over a 100-year climate run forced with $1W/m^2$ (left) and $0.5W/m^2$ (right) yearly increase in global heat flux forcing, showing the difference between the time-averaged last 5 years and first 5 years.
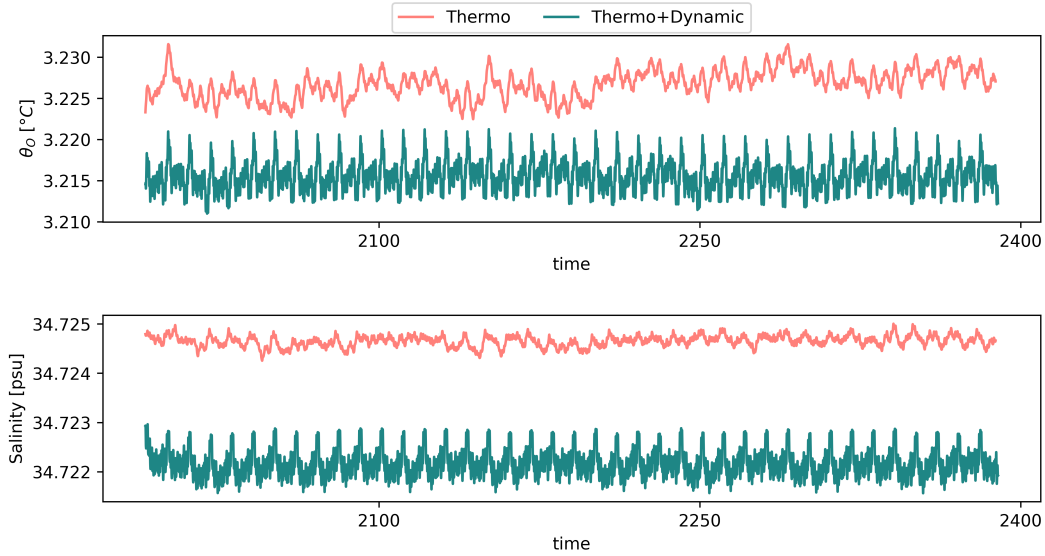
**Figure S15.** Spatially averaged potential temperature ($\theta_O$) and salinity ($S$) time series over a 400-year run forced with repeated atmospheric conditions taken from 1990-2000 for emulators $\mathcal{F}_{\text{thermo}}$ and $\mathcal{F}_{\text{thermo+dynamic}}$. The time series is averaged over 300-day rolling windows for visual clarity. The potential temperature trends for $\mathcal{F}_{\text{thermo}}$ and $\mathcal{F}_{\text{thermo+dynamic}}$ are $7.39 \times 10^{-6}$ $^\circ C$/year and $4.08 \times 10^{-7}$ $^\circ C$/year, respectively, while the Salinity trends are $1.867 \times 10^{-7}$ psu/year and $-1.397 \times 10^{-8}$ psu/year, respectively.
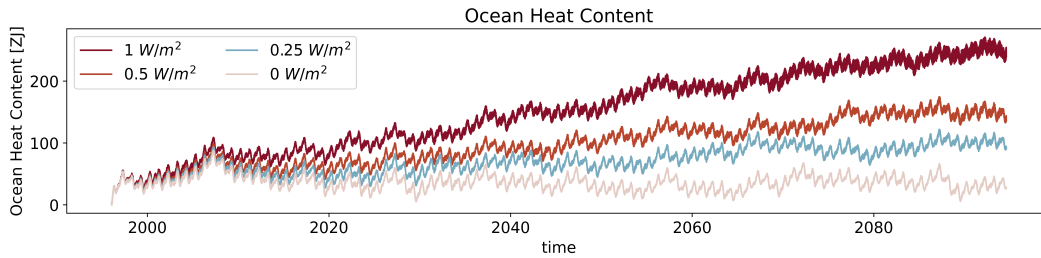


**Figure S16.** Ocean heat content trends for 100 year runs from the $\mathcal{F}_{\text{thermo}}$ emulator. These runs are forced by increasing the global heat flux, forcing $0, 0.25, 0.5$, and $1W/m^2$ per year to show how the emulator responds under a range of warming conditions.