

MIT Open Access Articles

*A Non#Intrusive Machine Learning Framework for
Debiasing Long#Time Coarse Resolution Climate
Simulations and Quantifying Rare Events Statistics*

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

Citation: Barthel Sorensen, B., Charalampopoulos, A., Zhang, S., Harrop, B. E., Leung, L. R., & Sapsis, T. P. (2024). A non-intrusive machine learning framework for debiasing long-time coarse resolution climate simulations and quantifying rare events statistics. *Journal of Advances in Modeling Earth Systems*, 16, e2023MS004122.

As Published: 10.1029/2023ms004122

Publisher: American Geophysical Union

Persistent URL: <https://hdl.handle.net/1721.1/154212>

Version: Final published version: final published article, as it appeared in a journal, conference proceedings, or other formally published context

Terms of use: Creative Commons Attribution-NonCommercial-NoDerivs License





RESEARCH ARTICLE

10.1029/2023MS004122

Key Points:

- Development of non-intrusive correction operators for coarse scale climate simulations
- Design of a training procedure that improves dynamics and allows for characterization of extremes with return period longer than the training data
- Application to Energy Exascale Earth System Model and demonstration of improvement on global and regional statistics

Correspondence to:

T. P. Sapsis,
sapsis@mit.edu

Citation:

Barthel Sorensen, B., Charalampopoulos, A., Zhang, S., Harrop, B. E., Leung, L. R., & Sapsis, T. P. (2024). A non-intrusive machine learning framework for debiasing long-time coarse resolution climate simulations and quantifying rare events statistics. *Journal of Advances in Modeling Earth Systems*, 16, e2023MS004122. <https://doi.org/10.1029/2023MS004122>

Received 14 NOV 2023

Accepted 25 FEB 2024

Author Contributions:

Conceptualization: L. R. Leung,

T. P. Sapsis

Data curation: S. Zhang**Formal analysis:** B. Barthel Sorensen,

S. Zhang, B. E. Harrop, L. R. Leung,

T. P. Sapsis

Funding acquisition: L. R. Leung,

T. P. Sapsis

Investigation: B. Barthel Sorensen,

A. Charalampopoulos, S. Zhang,

B. E. Harrop, L. R. Leung, T. P. Sapsis

Methodology: B. Barthel Sorensen,






A. Charalampopoulos, S. Zhang,

B. E. Harrop, L. R. Leung, T. P. Sapsis

© 2024 The Authors. Journal of Advances in Modeling Earth Systems published by Wiley Periodicals LLC on behalf of American Geophysical Union.

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial-NoDerivs License](#), which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

A Non-Intrusive Machine Learning Framework for Debiasing Long-Time Coarse Resolution Climate Simulations and Quantifying Rare Events Statistics

B. Barthel Sorensen¹ , A. Charalampopoulos¹, S. Zhang² , B. E. Harrop² , L. R. Leung² , and T. P. Sapsis¹ 

¹Department of Mechanical Engineering, Massachusetts Institute of Technology, Cambridge, MA, USA, ²Pacific Northwest National Laboratory, Richland, WA, USA

Abstract Due to the rapidly changing climate, the frequency and severity of extreme weather is expected to increase over the coming decades. As fully-resolved climate simulations remain computationally intractable, policy makers must rely on coarse-models to quantify risk for extremes. However, coarse models suffer from inherent bias due to the ignored “sub-grid” scales. We propose a framework to *non-intrusively* debias coarse-resolution climate predictions using neural-network (NN) correction operators. Previous efforts have attempted to train such operators using loss functions that match statistics. However, this approach falls short with events that have longer return period than that of the training data, since the reference statistics have not converged. Here, the scope is to formulate a learning method that allows for correction of dynamics and quantification of extreme events with longer return period than the training data. The key obstacle is the chaotic nature of the underlying dynamics. To overcome this challenge, we introduce a dynamical systems approach where the correction operator is trained using reference data and a coarse model simulation *nudged* toward that reference. The method is demonstrated on debiasing an under-resolved quasi-geostrophic model and the Energy Exascale Earth System Model (E3SM). For the former, our method enables the quantification of events that have return period two orders longer than the training data. For the latter, when trained on 8 years of ERA5 data, our approach is able to correct the coarse E3SM output to closely reflect the 36-year ERA5 statistics for all prognostic variables and significantly reduce their spatial biases.

Plain Language Summary We present a general framework to design machine learned correction operators to improve the predicted statistics of low-resolution climate simulations. We illustrate the approach, which acts on existing data in a post-processing manner, on a simplified prototype climate model as well as a realistic climate model, namely the Energy Exascale Earth System Model (E3SM) with 110 km resolution. For the latter, we show that the developed approach is able to correct the low-resolution E3SM output to closely reflect the climate statistics of historical observations as quantified by the ERA5 data set. We also demonstrate that our model significantly improves the prediction of atmospheric rivers, an example of extreme weather events resolvable by the low resolution model.

1. Introduction

As climate changes, several studies have indicated that the frequency and severity of extreme weather events will increase over the coming decades (Fischer et al., 2021; Raymond et al., 2020; Robinson et al., 2021). Accurately quantifying the risk of such events is a critical step in developing strategies to prepare for and mitigate their negative impacts on society—which can include billions of dollars in damages and thousands of lost lives (Allen et al., 2012; Fiedler et al., 2021; Houser et al., 2015). However, predicting the risk, magnitude, and impacts of such events is difficult and multifaceted. First, these events are seldom observed and arise due to a range of—often not fully understood—physical mechanisms (Lucarini et al., 2016; Sapsis, 2021). Moreover, the most devastating events are those which arise due to extreme excursions of multiple variables simultaneously, such as concurrent drought and heatwaves, which have a combined effect greater than each would have had in isolation (Bevacqua et al., 2023; Raymond et al., 2020; Robinson et al., 2021; Zscheischler et al., 2018). In addition, these extremes, whether occurring in isolation or in concert, interact with the earth system—and society—in myriad and often non-trivial ways. For example, the aforementioned combination of excess heat and below-average precipitation can increase the frequency of wildfires, degrade soil quality, and intensify water shortages, all of which then in turn have devastating socioeconomic impacts through, for example, reduced crop yields and even

Project administration: L. R. Leung, T. P. Sapsis
Resources: S. Zhang, L. R. Leung, T. P. Sapsis
Software: B. Barthel Sorensen, A. Charalampopoulos
Supervision: L. R. Leung, T. P. Sapsis
Validation: B. Barthel Sorensen, A. Charalampopoulos, S. Zhang, B. E. Harrop
Visualization: B. Barthel Sorensen, S. Zhang, B. E. Harrop
Writing – original draft: B. Barthel Sorensen, A. Charalampopoulos, S. Zhang, T. P. Sapsis
Writing – review & editing: B. E. Harrop, L. R. Leung

increased spread of disease (Barriopedro et al., 2011; Geirinhas et al., 2021; Hauser et al., 2016; Witte et al., 2011). Fully quantifying this complicated and interconnected system of physical, ecological, and social factors will surely require innovation and collaboration on a vast scale (Bauer et al., 2021; Slingo et al., 2022), yet even the first step, the accurate modeling of the climate dynamics, remains a challenging and unsolved problem.

At their heart, climate models (Manabe et al., 1965; Mintz, 1968; Smagorinsky, 1963; Smagorinsky et al., 1965), or their more modern counterpart, Earth System Models (ESM) (Dennis et al., 2012; Golaz et al., 2022; Taylor et al., 2009) are discretized forms of the equations of motion governing the Earth atmosphere and oceans. These known dynamical equations are then coupled to theoretical or empirical parameterizations of phenomena whose governing equations are unknown, such as the exact relationship between the vertical distribution of water vapor and precipitation rates (Holloway & Neelin, 2009; Stensrud, 2007) or the residence time of carbon in various terrestrial reservoirs (Bloom et al., 2016; Friend et al., 2014). Statistical climate predictions are then made by averaging over ensembles of realizations generated by such models. Unfortunately, a significant challenge in the practical application of these models is the computational complexity incurred by the vast range of dynamically active scales present in the oceans and atmosphere. This challenge is compounded when considering the need for large ensembles of models to be run over time horizons stretching decades or even centuries. The current state-of-the-art for climate modeling corresponds to an atmospheric spatial resolution of approximately 1° (i.e., approximately 110 km), with some early progress seen in the development of <5 km resolution models (Stevens et al., 2019; Tomita et al., 2005; Wedi et al., 2020). While there are some proponents of even finer (1 km) resolution simulations (Bauer et al., 2021; Slingo et al., 2022), even these fail to resolve critical phenomena such as the dynamics of stratocumulus clouds, which evolve on length scales of around 10 m (Schneider, Teixeira, et al., 2017; Wood, 2012), much less than the Kolmogorov dissipation scale which is on the order of 1 mm. In fact, the degrees of freedom in an ESM with 1 km resolution, which is stretching today's computational capabilities, fall short of what is needed to fully resolve atmospheric turbulence by a factor of 10¹⁷ (Schneider et al., 2023). These realities imply that the brute-force computation of the climate system will remain out of reach for the foreseeable future and that meaningful progress will require new and innovative solutions.

One promising and growing area of research to sidestep the computational intractability of fully resolved simulations is the combination of existing climate models with modern machine learning (ML) and data-assimilation strategies which learn the “sub-grid” dynamics from targeted high resolution simulations or observational data (Schneider et al., 2023; Schneider, Lan, et al., 2017). For example, reservoir-computing-based hybrid models have recently been demonstrated which learn online corrections to coarse climate models. These have been shown to substantially reduce overall bias (Arcomano et al., 2022) and capture events, such as sudden stratospheric warming, which are not resolved at all in free-running coarse climate models (Arcomano et al., 2023). Another, and perhaps more widely adopted approach is the data-driven parametric closure model. Here “closure model” refers to a state-dependent forcing term which aims to mimic the dynamic effects of the un-resolved scales on the resolved ones. Initially, such strategies were demonstrated on idealized aqua planet configurations using random forests (Yuval & O’Gorman, 2020) and neural network (NN) models (Brenowitz & Bretherton, 2019; Rasp et al., 2018; Yuval et al., 2021). More recently they have been applied to realistic global climate models to learn parametric forcing terms from reanalysis data using RFs (Watt-Meyer et al., 2021) and Deep Operator Networks (DeepONet) (Bora et al., 2023), as well as from higher resolution simulations with 3 km (Bretherton et al., 2022), and 25 km (Clark et al., 2022) resolution—both utilizing NNs and RFs. Across these studies, the ML closure models led to a robust improvement of 20%–30% in certain predicted integral quantities such as mean precipitation. However, predictions of other quantities were less reliable. For example (Clark et al., 2022), found that surface temperature predictions depended non-trivially on the random seed used in training the ML model. Furthermore, these approaches did not universally reduce the bias of the predicted climate over the uncorrected baseline, even in some cases increasing the bias of the coarse model (Clark et al., 2022; Watt-Meyer et al., 2021).

Despite these concerns, the most severe limitation of these approaches is numerical instability when integrating over long time horizons. This means that the aforementioned studies have only been demonstrated over short, 1 year (Watt-Meyer et al., 2021) and 5.25 year (Clark et al., 2022) time horizons—far shorter than what is required for long-term climate analysis. Such instabilities are inherent in this type of intrusive approach, except of special classes of representations for the closure terms which can guarantee stability of one-point and two-point statistics (H. Zhang et al., 2021). The ML correction term augmenting the coarse-scale equations is designed to bring the turbulent attractor of the corrected system in line with that of the reference. However, the ML approximation of the sub-grid scale dynamics will not be perfect, and due to the chaotic nature of the system, small excursions will

eventually grow, causing the predicted system trajectory to diverge from the attractor of the reference data (Wikner et al., 2022). We refer the interested reader to Yuval et al. (2021) for a detailed discussion of the stability challenges inherent in data-driven closure models.

Motivated by the intrinsic limitation of data-driven closure-models, we consider a different strategy. We seek to learn a ML operator which does not alter the equations, but rather acts as a post-processing operation to debias coarse scaled climate models. Such a *non-intrusive* approach has several theoretical advantages. First, it does not require altering the code of the core climate model—a non-trivial endeavor which often requires the harmonization of codes written in different languages (J. McGibbon et al., 2021). Second, unlike a closure model, it is domain agnostic, it can be applied globally or only for specific regions or altitudes. Third, and most critically, it is not susceptible to the same instabilities which plague schemes which apply ML corrections directly to the system dynamics. This in turn means it can be used to generate ensembles of trajectories over century + time horizons—a necessary step for quantifying risk of rare climate events with very long return periods. However, ML such a non-intrusive correction presents several considerable challenges, the foremost of which is the chaotic character of the climate systems under investigation. A mapping learned directly from some particular trajectory of a coarse model to a reference is unlikely to generalize, as it will encode not only the differences inherent in the coarse-scaling but it will also be corrupted by the particular chaotic realization of the training data. To overcome this challenge, Arbabi and Sapsis (2022) developed a generative framework which uses a system of linear stochastic differential equations in conjunction with a nonlinear map modeled through optimal transport. The nonlinear map and the stochastic linear system are optimized so that the statistics of the output match the statistics of the training data. In a more recent work, Blanchard et al. (2022) used a more complex architecture consisting of a spatial wavelet decomposition, a temporal-convolutional-network (TCN) and long-short-term-memory (LSTM) architectures trained also on a purely statistical loss function involving single point probability densities and temporal spectrum. Alternatively, strategies such as generative adversarial networks (J. J. McGibbon et al., 2023) and unsupervised image-to-image networks (UNIT) (Fulton et al., 2023) have been used to correct biases in average precipitation rates—an integral quantity which is less affected by stochastic variation. While ML correction operators using a purely statistical loss function can indeed generate trajectories with plausible statistics, this property alone does not guarantee the resulted spatio-temporal dynamics are always physically realistic. Most importantly the quality of the resulted models, by design, cannot exceed the quality of the statistics used for training. Therefore, if the statistics for rare events of a given (large) return period have not converged (because of low availability of such events in the training set) the model is essentially forced to reproduce inaccurate, that is, non-converged statistics, at least for rare events that have return periods comparable or longer than the training data set. To this end, methods based on purely statistical loss functions cannot be used for statistical extrapolation.

In this work we describe a framework to overcome this challenge. Our aim is to design an algorithm that learns essential dynamics and is able to extrapolate statistics with a non-intrusive approach. The heart of the proposed strategy is that we do not machine learn a map from any *arbitrary* coarse trajectory to the reference, but specifically from a coarse trajectory *nudged towards that reference*. Nudging the coarse model towards the target reference trajectory results in an input trajectory which predominately obeys the dynamics of the coarse model yet remains close to the reference trajectory. Training a ML operator on this specific pair of trajectories allows us to learn a transformation which encodes only the differences caused by the coarse-grid without being corrupted by random stochastic effects. Once trained, this correction operator can then reliably map *any* free-running coarse trajectory into the attractor of the reference data. We first lay out the theoretical framework of the proposed strategy in terms of a general chaotic dynamical system. We then illustrate our method on a simplified 2-layer quasi-geostrophic (QG) model, and show that we are able to correct a severely under-resolved solution to accurately reflect the long time statistics of the fully resolved reference—even when the model is trained on much shorter time histories than the reference. Finally, we apply our framework to a realistic climate model, the Energy Exascale Earth System Model (E3SM) with ~ 110 km grid resolution. We show that using only 8 years of training data our correction operator is able to bring the global and regional 30-year statistics of the primitive variables into good agreement with ERA5 reanalysis data, and reduce the error in the 36-year average integrated vapor transport (IVT) by 51% relative to the free-running E3SM solution. Our results show that our framework is able to characterize statistics of events with a return period that is multiple times longer than the length of the training data and therefore represent a promising step towards reliable long term climate predictions.

The remainder of the article is organized as follows. In Section 2 we introduce the mathematical framework and general ML strategy. We then apply our method to a quasi-geostrophic model in Section 3 and the E3SM climate

model in Section 4. Finally we conclude with a discussion of the implications of our results and the potential extensions and limitations of our method in Section 5.

2. Training Correction Operators for Imperfect Chaotic Systems

We consider a high-resolution discretization of an ergodic chaotic dynamical system, and its solution (named thereafter the reference solution),

$$\dot{\mathbf{u}} = F(\mathbf{u}), \quad \mathbf{u} \in \mathbb{R}^N \quad (1)$$

as well as, a coarse discretization of the same dynamical system (referred as CR), described by the model

$$\dot{v} = f(v), \quad v \in \mathbb{R}^n, \quad (2)$$

where $n < N$. The reference solution is projected to the coarse grid through the projection operator \mathcal{P} , that is,

$$u = \mathcal{P}\mathbf{u}, \quad u \in \mathbb{R}^n \quad (3)$$

The objective of this work is to capture the long time statistics of u by solving the imperfect model (Equation 2) and then applying a correction operator, \mathcal{G} , to the computed solution. The correction operator is assumed to be spatially non-local, with memory, but causal, that is, the correction at time t may depend only on the past of the input but not the future. To learn this correction operator we assume a reference data set (referred as RD) generated by the high resolution model or reanalysis data in the form of a finite time trajectory: $\{u(t), t \in [0, T]\}$.

This is a non-trivial problem since any CR trajectory (Equation 2) and RD trajectory (reference data set U) will not be comparable, that is, cannot be used to formulate the training of the correction operator as a supervised learning problem. In fact, even if the initial condition of the imperfect model is chosen to be the same with u ($t = 0$), the two trajectories will rapidly diverge due to the chaotic nature of the system.

In Blanchard et al. (2022) the authors aim to address this fundamental obstacle by developing a cost function that penalizes directly the deviation between the generated statistics of $\mathcal{G}(v)$ and the statistics of the reference trajectory, u . While the approach has shown some promise, it is a very hard optimization problem that often results in non-physical realizations, $\mathcal{G}(v)$. At a more fundamental level, the approach does not really utilize the “sequencing” or dynamics encoded in the reference data, but rather its statistics, which for real world problems cannot be guaranteed to be accurate especially for rare events (e.g., using 40 years of reanalysis data cannot guarantee accurate statistics for rare events with a longer return period).

Here we follow a radically different method that aims to learn the correction operator \mathcal{G} using the reference trajectory and the dynamics of the coarse model, rather than their corresponding finite-time statistics. One of the key objectives of this work is the identification of a data set which will allow for the training of such a correction operator. The primary challenge therein is the need to suppress the chaotic divergence of the coarse scale model during the training phase.

We consider the deviation of the two dynamical systems:

$$q \equiv v - u, \quad q \in \mathbb{R}^n. \quad (4)$$

By computing the derivative we have an equation along the reference trajectory, \mathbf{u} ,

$$\dot{q} = f(v) - \mathcal{P}F(\mathbf{u}) = f(q + \mathcal{P}\mathbf{u}) - \mathcal{P}F(\mathbf{u}). \quad (5)$$

The right hand side expresses, for a given \mathbf{u} , the way the two models diverge. Naturally, the above equation will provide useful information between the two trajectories for as long these remain close to each other. Beyond that point, that is, after chaotic divergence has occurred, it is not meaningful to compare the two trajectories. To address this issue, we add a damping term in the right hand side of Equation 5 that will keep the deviation small:

$$\dot{q}_\tau = f(q_\tau + \mathcal{P}\mathbf{u}) - \mathcal{P}F(\mathbf{u}) - \frac{1}{\tau}q_\tau, \quad (6)$$

where τ is a constant relaxation time scale that is chosen so that the added term is at least one order of magnitude smaller compared with all the other terms in Equation 6. Moreover, we add the subscript τ to emphasize that this is divergence computed with the artificial damping term. The added term is large enough to guarantee that over time scales longer than τ the deviation does not grow exponentially due to chaotic effects, that is, the coarse scale model remains in a relevant state to the reference state, but also small enough to allow for the coarse scale model dynamics to evolve unimpeded. The last point is essential in order to obtain a data set with sufficient content regarding the imperfection of the coarse scale model.

By transforming the equation for q_τ into the v variable, we obtain the final equation for the generation of *nudged* data sets to be used for training:

$$\dot{v}_\tau = f(v_\tau) - \frac{1}{\tau}(v_\tau - u), \quad (7)$$

where the second term on the right hand side is known as the nudging tendency. The pair of trajectories (v_τ, u) is the basis for training the correction operator. We note that nudging has been widely used in the context of data-assimilation to improve the predictive capabilities of climate models (Huang et al., 2021; Miguez-Macho et al., 2005; Storch et al., 2000; Sun et al., 2019) as well as on developing hybrid approaches for climate modeling (Bretherton et al., 2022). Here the use of nudging is only for the development of relevant training pairs of trajectories.

2.1. Interpretation of training with data from the nudged model

To obtain a dynamical understanding of the mapping process between the nudged trajectory generated by the above equation and the exact trajectory, we hypothesize the existence of a slow-fast decomposition for v_τ and u . Our motivation is the observation that for many turbulent systems, spatially-coarse modeling affects primarily the fast time scales while it results in smaller errors in the slow time scales. However, fast time scales are important for the characterization of extreme events, as the latter are typically short lived structures. We express the solution v_τ in the following slow-fast decomposition based on the relaxation time scale τ :

$$v_\tau(t) = v_s(\mathcal{T}) + v_f(t), \quad (8)$$

where $\mathcal{T} = \epsilon t$ is the slow time scale, and $\epsilon = 1/\tau \ll 1$, where τ is the relaxation time scale. Moreover, we also decompose the reference solution in a slow-fast form:

$$u(t) = u_s(\mathcal{T}) + u_f(t), \quad (9)$$

Based on the above, we have by direct calculation:

$$\dot{v}_\tau(t) = \epsilon v'_s(\mathcal{T}) + \dot{v}_f(t, v_s), \quad \text{where } v' = \frac{dv}{d\mathcal{T}}. \quad (10)$$

Substituting into Equation 7 we obtain

$$\epsilon v'_s + \dot{v}_f = f(v_s + v_f) + \epsilon(u_s + u_f - v_s - v_f). \quad (11)$$

Separating the slowly evolving terms of order $\mathcal{O}(\epsilon)$, that is, the small terms that depend only on \mathcal{T} , we have:

$$v'_s = u_s - v_s \Rightarrow v_s(\mathcal{T}) = \int e^{-(\mathcal{T}-s)} u_s(s) ds. \quad (12)$$

The fast terms on the other hand will give, to zero order:

$$\dot{v}_f = f(v_s(\mathcal{T}) + v_f) + \mathcal{O}(\epsilon). \quad (13)$$

From the last two equations we can conclude that Equation 7 essentially drives the coarse scale model along the slow dynamics of the reference attractor captured by the trajectory, u , Equation 12, but leaves the fast dynamics free to evolve according to Equation 13. By driving the imperfect model in regions of the attractor where we have reference data we are able to define a supervised learning problem, where the input is the solution with imperfect fast dynamics defined by Equation 7 and the output is the reference solution, u . In this way, one can use this pair of data to machine learn a map that corrects the fast features of the imperfect model, where the largest model errors are concentrated due to coarse discretization.

It is important to emphasize that the method does not assume any scale separation in the dynamics. Instead the parameter τ controls which temporal scales are corrected by the NN operator. On the other hand, it is important to mention that the success of the scheme relies on a minimum data requirement, sufficient to guarantee proper generalization of the correction operator.

2.2. Selection of the relaxation time scale τ

One of the key questions in the practical implementation of this framework is the choice of the relaxation timescale τ . It quantifies the strength of the nudging tendency and represents a trade off between the suppression of the chaotic divergence and the suppression of the inherent dynamics of the coarse model. If $\tau \rightarrow \infty$, the nudging tendency in Equation 7 will be too weak to suppress the chaotic divergence between v_τ and u . This will mean that a map between them will not generalize when applied to free-running coarse solutions. Alternatively, if $\tau \rightarrow 0$, the nudging tendency will completely suppress the dynamics and v_τ will be indistinguishable from u and a map between them will be trivial. From numerical experiments we performed, we found that a value of τ that results in a nudging term that is one order of magnitude smaller than the other terms of the model represents a good balance between these extremes, that is, the performance of the algorithm remains the same as long as the choice of τ remains within this range.

2.3. Spectrum-matched nudging

Before we proceed to the ML of the correction operator we need to address an energetic inconsistency created by the inclusion of the nudging term in the coarse scale model. This is associated with the artificial dissipation that is introduced to the dynamics of the model due to the term $\frac{1}{\tau}v_\tau$. While the term is generally smaller than all other terms of the model, it still creates small discrepancies between the spectra of the nudged solution, v_τ , and the free coarse solution, v . This is an inconsistency that has been observed in different settings of data-assimilation and several solutions have been proposed, for example, 4DVar (Mons et al., 2016) or ensemble variational method (Buchta & Zaki, 2021; Mons et al., 2016).

Here we employ the simplest approach to correct the spectral inconsistency: we rescale the spectrum of the nudged trajectory, v_τ to match the spectrum of the coarse model spectrum. Specifically, let $\hat{u}_k = \mathcal{F}[u]$ be the spatial Fourier transform of the field u . We define the spectral energy as

$$\mathcal{E}_{k,u} = \frac{1}{T} \int_0^T |\hat{u}_k|^2 dt. \quad (14)$$

Next, we consider the energy-ratio per wavenumber, between the free-running, v , and the nudged solution, v_τ , defined as

$$a_k \equiv \sqrt{\frac{\mathcal{E}_{k,v}}{\mathcal{E}_{k,v_\tau}}} \quad (15)$$

We define as the spectrum-matched nudged solution as the inverse Fourier transform of the spectrally rescaled nudged solution:

$$v'_\tau = \mathcal{F}^{-1}[a_k \hat{v}_{k,\tau}]. \quad (16)$$

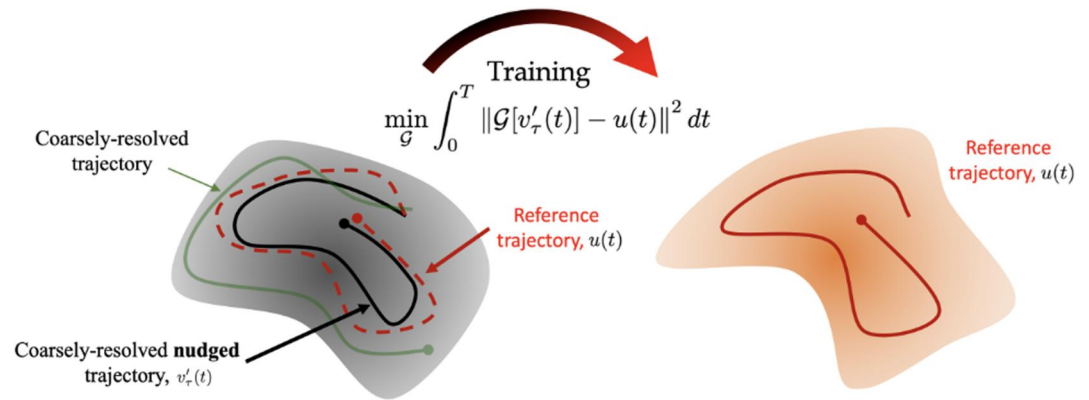


Figure 1. Description of the method that learns a map between the attractor of the coarsely-resolved equations and the attractor of the reference trajectory. Left: the red dashed curve represents the reference trajectory. The black curve is a coarsely-resolved nudged trajectory toward the reference trajectory. The green curve is the free-run coarsely-resolved trajectory that is not used for training (shown for reference). Right: the target attractor and the target trajectory (red), same as the dashed curve shown at the left plot.

The resulted pair of *spectrally-corrected nudged* solution, v'_τ referred in what follows as nudged coarse (NC) data set, together with the reference data set (RD), u define a supervised learning problem with cost function being:

$$\min_G \int_0^T \|\mathcal{G}[v'_\tau(t)] - u(t)\|^2 dt \quad (17)$$

The training framework is graphically illustrated in Figure 1. In contrast to previous approaches that aim to match the statistics of the transformed output with statistics of a reference trajectory, the above optimization problem encodes directly the dynamics that is, the time sequencing of the data set. This property is crucial for better generalization capabilities, that is, to train with a short data set and be able to capture statistics that correspond to much longer simulations.

After we have machine learned the correction operator, \mathcal{G} , we apply it to the free running coarse model trajectory (CR), $v(t)$. The result is then used to compute statistics and other properties of interest. The workflows for training and testing are summarized in Figure 2. We emphasize that nudging and reference data are used only in the training phase. At the testing phase, the model is using only free-running coarse data and transform it to obtain the correct statistics. The good generalization capabilities of the correction operator allow for its application on much longer time series than those used for training, that is, the characterization of extreme events with return period that is longer than the training data set.

3. Quasi-Geostrophic Model

3.1. Background

As a first example we apply the presented correction method to the two-layer incompressible quasi-geostrophic (QG) flow (Qi & Majda, 2018). In a dimensionless form, its evolution equation is given by

$$\frac{\partial q_j}{\partial t} + \mathbf{u}_j \cdot \nabla q_j + (\beta + k_d^2 U_j) \frac{\partial \psi_j}{\partial x} = -\delta_{2,j} r \nabla^2 \psi_j - \nu \nabla^8 q_j \quad (18)$$

where $j = 1, 2$ corresponds to the upper and lower layer respectively, r the bottom-drag coefficient and β is the beta-plane approximation parameter, and k_d^2 represents the deformation frequency which for this study we fix at 4—a value consistent with the radius and rotation of the earth and the characteristic length and velocity scales of the atmosphere (Qi & Majda, 2018). This model is intended to approximate mid to high latitude atmospheric flows subject to an imposed shear current. A Taylor expansion of the Coriolis force reveals that for this assumption to hold we require roughly that $\beta \in [1, 2]$, which corresponds to an approximate latitude range of $\phi_0 \in [29^\circ, 64^\circ]$.

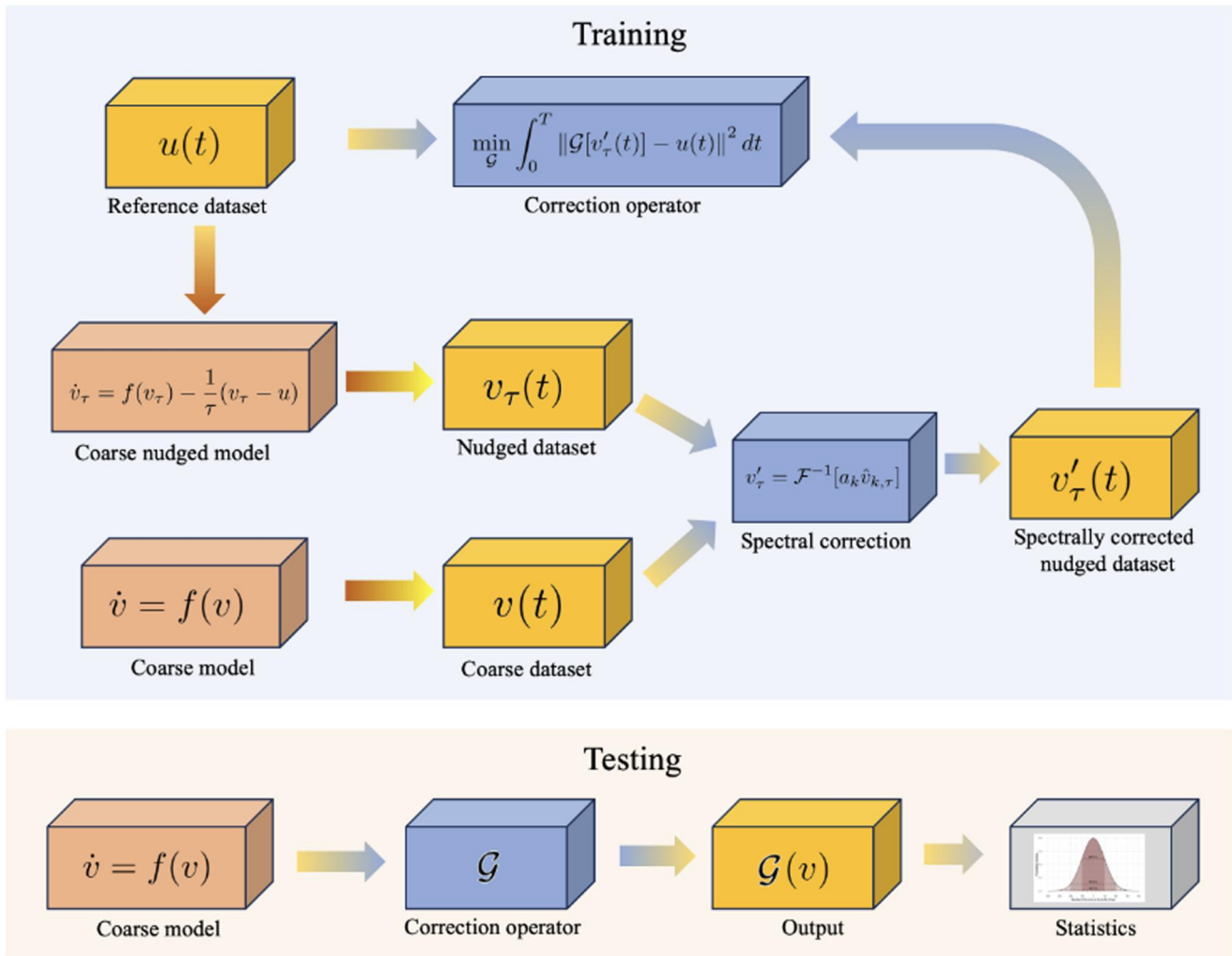


Figure 2. Workflow of the training process (top) and testing process (bottom) for the machine learning of correction operators and their application on the generation of long time climate simulations, that is, longer than the reference data set.

The flow is defined in the horizontal domain $(x, y) \in [0, 2\pi]$ and is subject to doubly periodic boundary conditions. The state variable is represented in three forms: velocity: \mathbf{u}_j , potential vorticity (PV): q_j and the stream function: ψ_j . The latter are related via the inversion formula

$$q_j = \nabla^2 \psi_j + \frac{k_d^2}{2} (\psi_{3-j} - \psi_j) \quad (19)$$

and the velocity is related to the stream function by $\mathbf{u}_j = U_j + \hat{\mathbf{k}} \times \nabla \psi_j$, where $\hat{\mathbf{k}}$ is the unit vector orthogonal to the (x, y) plane and $U_j = -1^{(j+1)} U$, with $U = 0.2$ represents the imposed mean shear flow. The corresponding nudged system of equations is given by

$$\frac{\partial q_j}{\partial t} + \mathbf{u}_j \cdot \nabla q_j + (\beta + k_d^2 U_j) \frac{\partial \psi_j}{\partial x} = -\delta_{2,j} r \nabla^2 \psi_j - \nu \nabla^8 q_j - \frac{1}{\tau} (q_j - q_j^{RD}) \quad (20)$$

where q_j^{RD} is the reference solution projected to the grid of q . We fix the nudging parameter $\tau = 16$ —a value for which we found the nudged solution tracks the reference, but generally retains the spectral properties of the free-running coarse solution. Furthermore, we note that while the nudging penalty is applied to the vorticity, it could

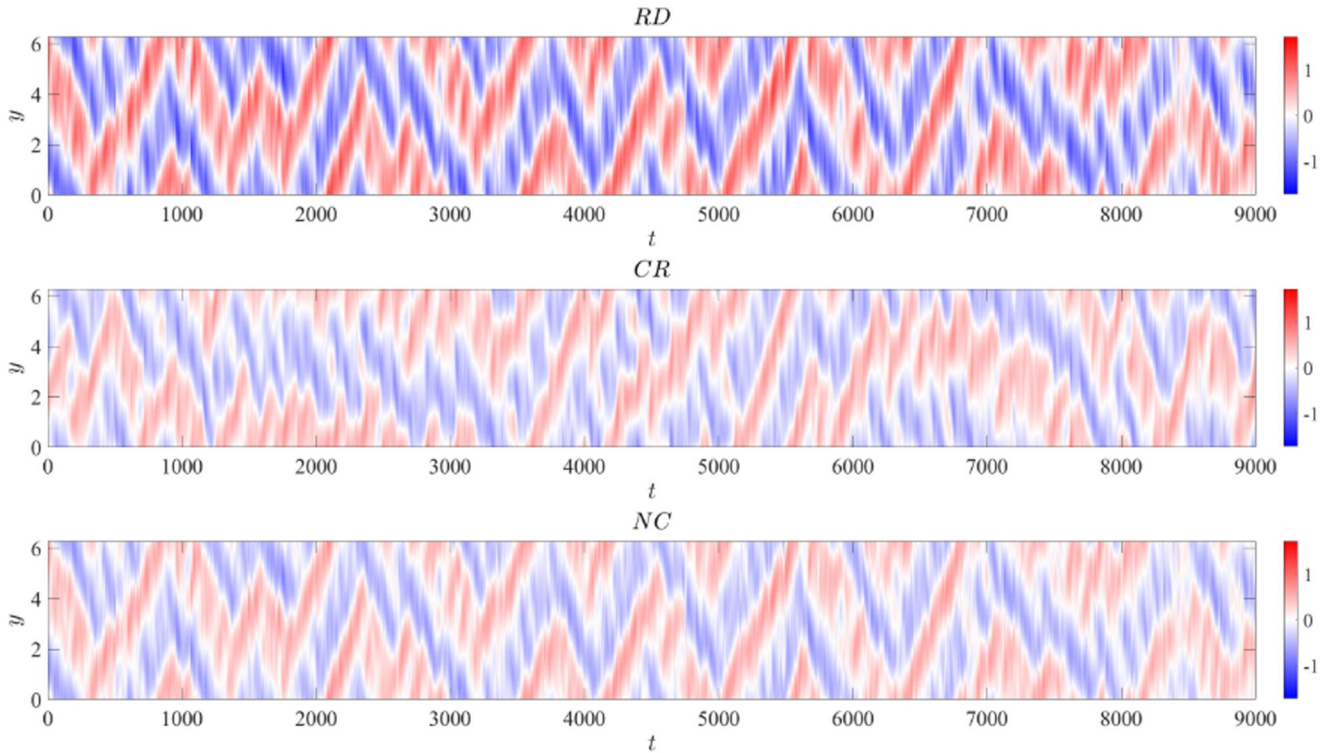


Figure 3. Example zonally averaged stream function $\bar{\psi}_1$ of the QG system Equation 18 for $\beta = 2.0$ and $r = 0.1$. From top to bottom: fully resolved, reference solution (RD), free-running coarse simulation (CR), and spectrally corrected nudged simulation (NC).

have equivalently been applied to the stream function or velocity. These possibilities are not explored in this work, however, as these three variables are all directly related we would not expect significant differences in the results.

Equations 18 and 20 are solved using a spectral method, with a spectral resolution of 24×24 and 128×128 for the coarse-scale (CR) and reference (RD) data respectively. The time integration is evaluated using a fourth order Runge-Kutta scheme with the same temporal resolution used for both the under- and fully-resolved simulations. Throughout the following discussion all results will be presented in the form of the stream function—as this uniquely defines the velocity and thus vorticity, this choice incurs no loss of generality. Additionally, we define the zonally averaged stream function as the integral over the x dimension,

$$\bar{\psi}_j(y, t) = \frac{1}{2\pi} \int_0^{2\pi} \psi_j(x, y, t) dx. \quad (21)$$

In Figure 3 we show the zonally averaged stream function in layer 1 for $\beta = 2.0$ and $r = 0.1$ of the three data sets: RD, CR, NC, as an illustrative example of both the fully- and under-resolved solutions. The primary qualitative difference between the coarse and fine grid solutions is in their amplitude. This is particularly clear when comparing the tails of the distributions in 3b. Note that the spectrally corrected NC solution reflects the qualitative spatio-temporal behavior of the fully resolved (RD) solution but exhibits the lower magnitude of the coarse (CR) solution.

3.2. NN Architecture and Training Strategy

The NN model we employ as a correction operator takes as an input the stream function field of both layers which is of dimension $24 \times 24 \times 2$. This vector is then compressed through a fully connected layer of dimension 60 and then passed through a long-short-term-memory (LSTM) layer of the same size before being expanded through a second fully connected layer to restore the data to its original size. The fully connected layers utilize hyperbolic tangent activation and the LSTM layer uses a hard-sigmoid activation. The model is trained purely on stream function data and thus the output of the model represents the statistically corrected stream function field.

The model is trained on a semi-physics informed loss function which consists of the L^2 norm of the error augmented with a second term which penalizes errors in the conservation of mass.

$$L = \sum_{j=1}^2 \int_0^{2\pi} \int_0^{2\pi} |\psi_j^{ml} - \psi_j^{rd}|^2 dx dy + \sum_{j=1}^2 \int_0^{2\pi} \int_0^{2\pi} \psi_j^{ml} dx dy \quad (22)$$

Here ψ^{ml} and ψ^{rd} denote the machine learned prediction (i.e., the ML transformation of the nudged data set) and the reference stream functions respectively. The mass conservation term is derived by noting that the two stream functions are linearly related to the height disturbances of the two layers and that by conservation of volume the integral of all height disturbances must vanish.

The correction operator is trained for 2,000 epochs on sequences of 100 data points spanning 10 time units taken from a single realization of the flow with $\beta = 2.0$ and $r = 0.1$ of length 1,000 time units. We then apply the trained correction operator to a separate (unseen) realization of the flow to generate the following results. .

3.3. Results

3.3.1. Prediction of Long Time Statistics

First, we apply our models, which are trained on data with $\beta = 2.0$ and $r = 0.1$, to a new realization of the flow with these same parameters. *A key objective of this work is to compute extreme event statistics for events that have a return period that is longer than the length of the training data.* Therefore, the question is how accurately we can capture the tails with a corrected long realization of the coarse model, when the correction operator has been trained on data that does not accurately the tails, that is, data of limited length.

To this end, we first apply our ML correction operator, which is trained on $T_{train} = 1,000$ time units of data, to a new realization of the flow spanning $T_{test} = 34,000$ time units. Figure 4a shows the global power spectra and probability density functions of the stream function in both layers. The power spectra are computed by taking the spatial average of the point-wise temporal power spectra, and the probability density function is taken across all space and time. The fully-resolved (RD) and under-resolved (CR) solutions are shown in solid and dashed black respectively and the ML correction of the under-resolved solution, henceforth denoted ML(CR), is shown in blue. As a reference, we also plot the statistics of the training data (RD_{train}) in red.

For both layers, the ML correction brings the coarse solution into good agreement with the fully-resolved reference. In terms of the spectra, the ML correction accurately captures the two peaks around $f = 0.15$, and only deviates significantly at very high frequencies. In terms of the probability density functions, the model slightly underpredicts the positive tail in layer 2, but captures the general shape well. Crucially, we note that the statistics of the (1,000 time unit) training data are meaningfully different from the (34,000 time unit) test data used to generate the results. Note especially the severe under-resolution of the spectrum and the discrepancy of the far tails of the probability density functions. This highlights the capability of our approach to capture tail events which are not present in the training data, most notably in layer 1. This is an important feature, as any practical long term (100+ year) climate prediction will necessarily be trained on far less training data. Furthermore, this highlights the advantages of our approach to one such as (Blanchard et al., 2022) in which the ML correction operator is trained to purely reproduce statistics, as such an approach is by construction restricted to the statistics of the training data.

Beyond capturing the global statistics, it is crucial for our model to accurately capture the dynamics evolving at specific spatial scales. Therefore, in Figure 5 we show the probability density function of a selection of the individual Fourier modes, parameterized by the wavenumber vector $\mathbf{k} = [k_x, k_y]$. In the interest of space we show the probability density of the barotropic stream function, defined as the average of the two layers. In general, the model captures the probability distributions of the Fourier modes very well, with some discrepancy in the far tails. Interestingly, the ML correction tends to underestimate the tails of the largest modes for example, $\mathbf{k} = [0, 1]$, and $[1, 0]$, while then trending toward overestimating the tails of the smaller modes for example, $\mathbf{k} = [2, 1]$, and $[2, 2]$.

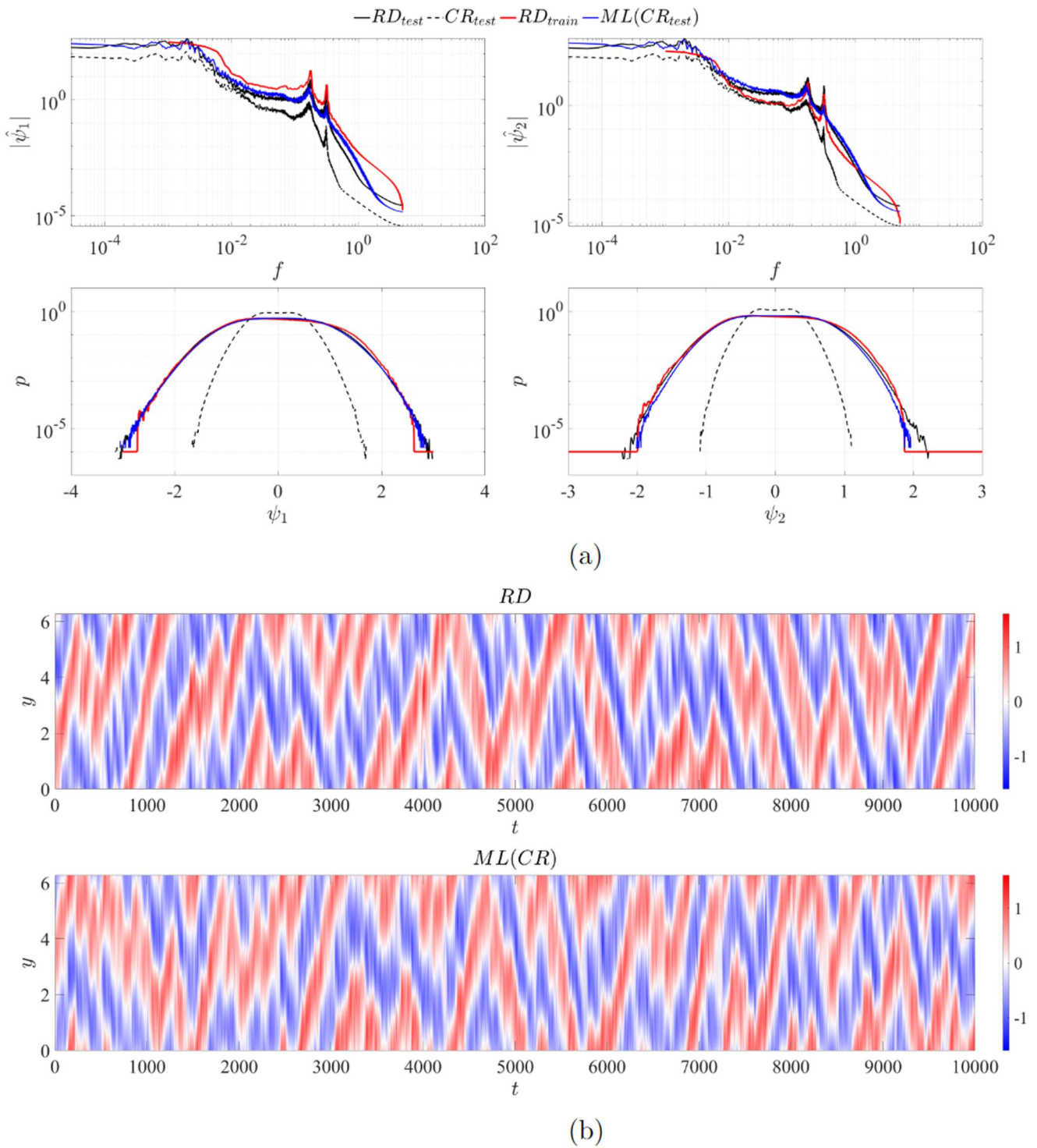


Figure 4. Model prediction for $\beta = 2.0$ and $r = 0.1$. Power spectrum and probability density function of stream functions ψ_1 (top row) and ψ_2 (bottom row). Test data, RD (solid black), CR (dash black), ML(CR) (blue) and training data RD_{train} (red) (a). Zonally averaged stream function $\bar{\psi}_1$, RD (upper panel) and ML(CR) (lower panel) (b). $T_{train} = 1,000$ and $T_{test} = 34,000$.

Finally, we reiterate that the only claim we make upon the trajectories predicted by our model is that they reflect the statistical properties of the fully resolved system. However, we expect our predictions to exhibit the qualitative behavior of the exact solution. To this end we show in Figure 4b a 10,000 time unit interval of the zonal average of

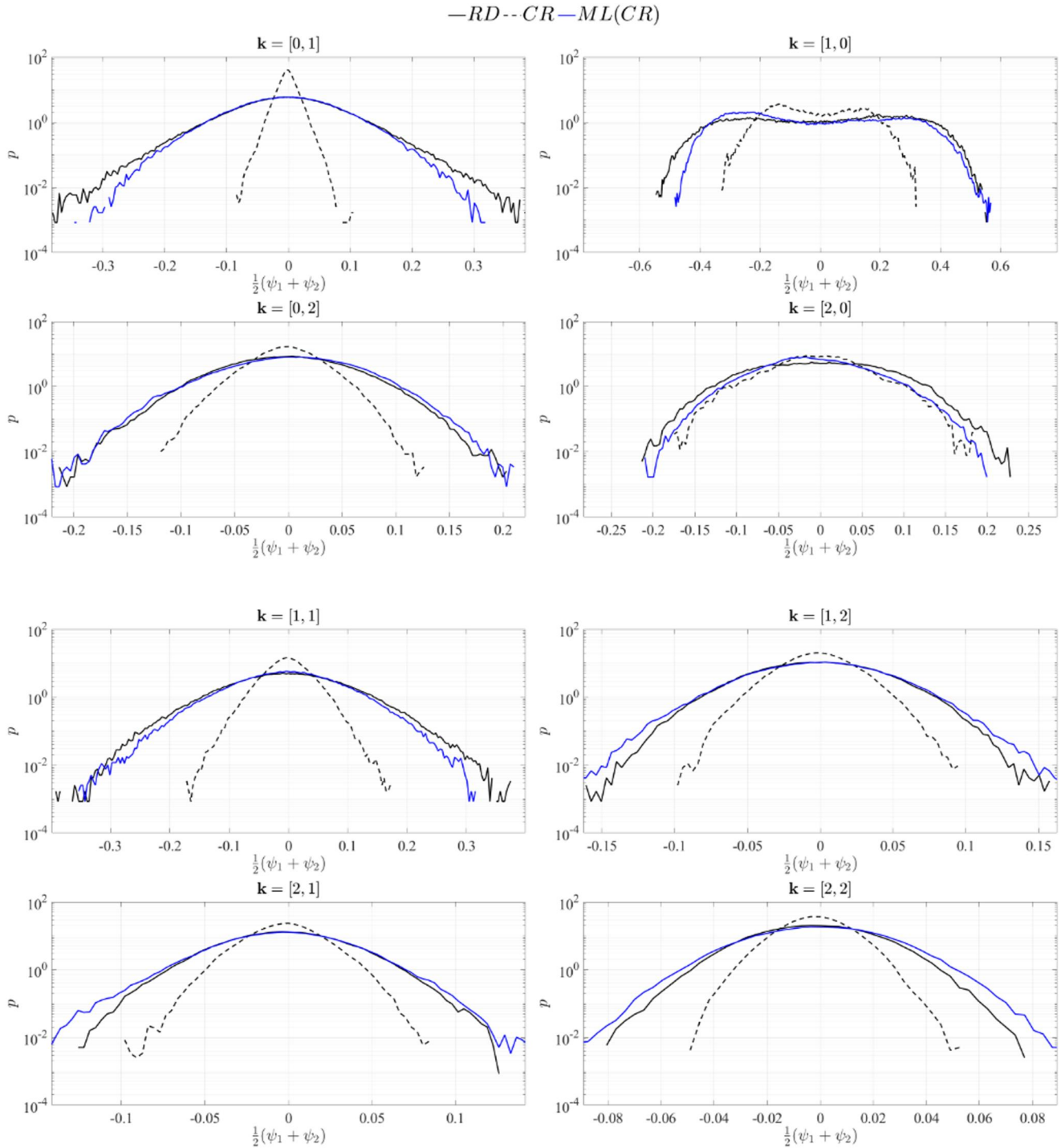


Figure 5. Probability density function of individual Fourier modes for $\beta = 2.0$ and $r = 0.1$. RD (solid black), CR (dashed black), ML(CR) (blue). $T_{train} = 1,000$ and $T_{test} = 34,000$.

the predicted solution. We do not show the full 34,000 time unit time horizon in order to improve the readability of the figure and highlight the spatiotemporal structure of the flow. We do indeed find good qualitative agreement with the fully-resolved simulation across the full test trajectory.

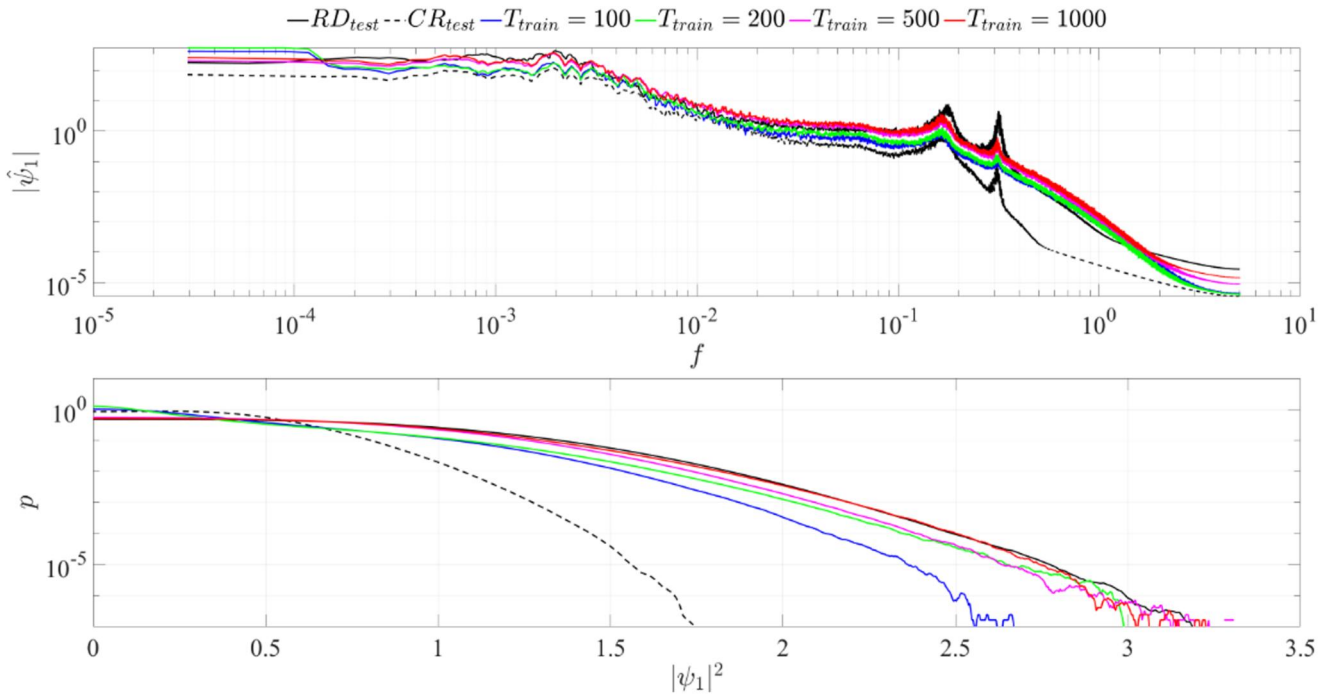


Figure 6. Model prediction of power spectrum and probability density function of $|\psi_1|$ for $T_{train} = 100, 200, 500,$ and $1,000$. For all cases $T_{test} = 34,000$.

3.3.2. Minimum Training Data Requirement

In the previous section we showed that our ML operator is capable of correcting the tails of a long time horizon coarse solution even when trained on a far shorter span of data. Here we investigate the minimum amount of training data needed to capture the long time ($T_{test} = 34,000$ time unit) statistics. We compare the results of our ML correction operator trained on data spanning $T_{train} = 100, 200, 500,$ and $1,000$ time units—the latter corresponding to the results described above. Both training and testing is carried out on data with $\beta = 2.0$ and $r = 0.1$. The probability density function and power spectrum of $|\psi_1|$ for these four cases are shown in Figure 6. We focus on the probability density function of the absolute value of the stream function in the interest of brevity. We see that the ML operator requires a minimum T_{train} between 500 and 1,000. While, the ML operators trained on $T_{train} < 500$ do improve the statistics of the coarse model, they do not capture the tails of the pdf and also underpredict the two spectral peaks. This is consistent with a closer examination of Figure 3 which shows that the characteristic time scale over which the large scale motions of the flow evolve is approximately 500–1,000 time units. Thus, for the QG model considered here, the ML operator requires seeing at least one full characteristic period of the flow in training. However, once it has seen one or two it is capable of learning the general features of the flow and can accurately reproduce statistics over much longer time horizons. This is a critical observation since for climate models data is always limited in time and the existence of such critical threshold can indeed pave the way for the computation of statistics for events that have return period much longer than the training data.

3.3.3. Evaluation for Different Flow Parameters Than the Training Data

Next, we apply the same ML operator to a realization of the QG model with flow parameters which differ from the training data, namely $\beta = 1.1$ and $r = 0.5$. For these parameter choices the flow lacks the characteristic spectral peaks of the β and r_d used to train the model exhibiting much more uniform frequency content. The lack of a dominant (slower) time scale means the flow evolves on faster characteristic time scale than the training data. These features make this a challenging test case to evaluate the generalizability of our model. Due to the shorter characteristic time scales, and the associated increased computational cost, for this experiment we consider a test data set of length $T_{test} = 10,000$ time units.

The results are summarized in Figures 7 and 8. In the former we plot the power spectra and probability density function and in the latter we plot the scale-by-scale probability density functions. In terms of the global statistics,

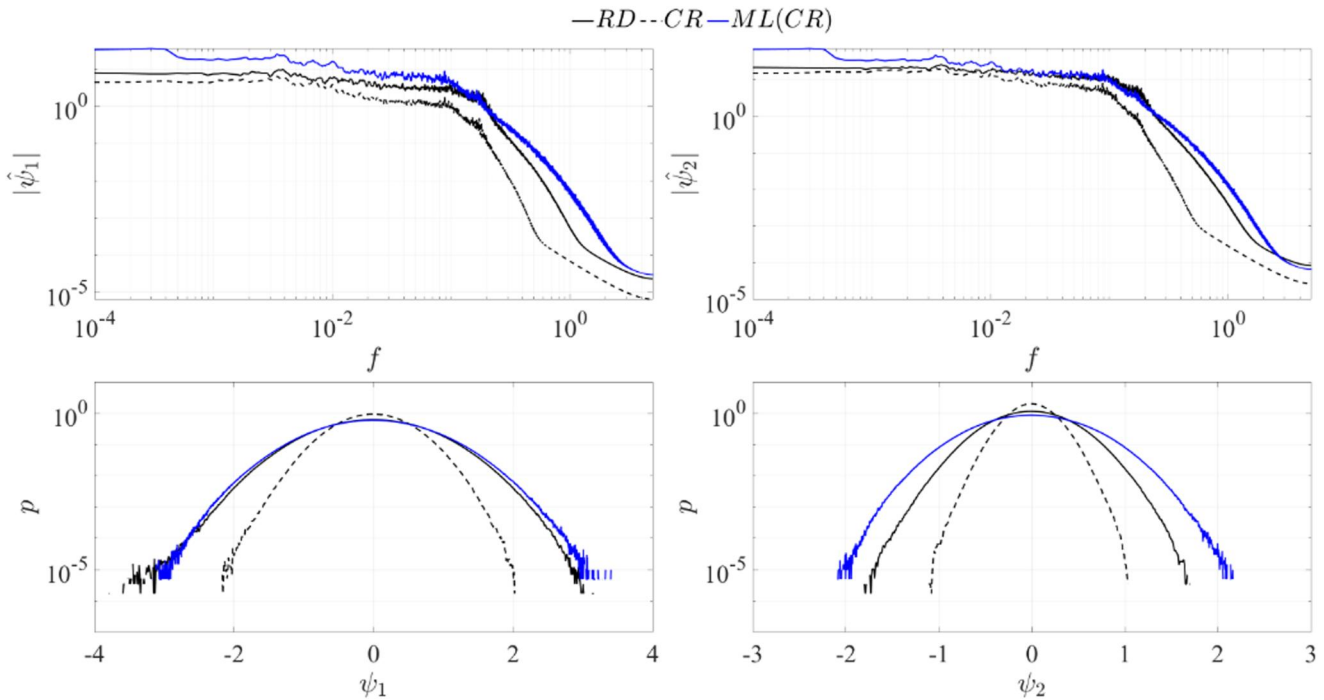


Figure 7. Model prediction for $\beta = 1.1$ and $r = 0.5$. Power spectrum and probability density function of stream functions ψ_1 (left) and ψ_2 (right), RD (solid black), CR (dash black), ML(CR) (blue). Training data: $\beta = 2.0$ and $r = 0.1$.

the predicted spectrum is in good agreement with the reference across much of the frequency domain, but underpredicts the spectral decay, and thus over-predicts the strength of the highest frequencies. In terms of the probability density function, there is excellent agreement in layer 1, while in layer 2 the model notably over-predicts the tails. The predictions of the scale-by-scale statistics are reasonably accurate and provide significant improvement over the free-running coarse model. However, the ML correction tends to over emphasize the strength of the tails for the larger length scales, for example, $\mathbf{k} = [0, 1], [1, 0], [1, 1]$. This is not surprising finding given the drastic over-correction of the tails in layer 2 seen in Figure 7.

4. Global Climate Model

4.1. Data Set

We now apply our framework to a realistic global climate model, the Energy Exascale Earth System Model (E3SM). In particular, version 2 of the E3SM Atmosphere Model (EAMv2) (Dennis et al., 2012; Golaz et al., 2022; Taylor et al., 2009). The progress variable is $\mathbf{X}(\theta, \phi, k, t) = (U, V, T, Q)$. The progress variables (U, V) correspond to the zonal and meridional components of wind velocity, T is air temperature and Q is specific humidity. The spatial coordinates (θ, ϕ, k) are the polar, $\theta \in [-90, 90]$, azimuthal angles, $\phi \in [0, 360]$, and the sigma level respectively. The latter of which can be understood as a measure of altitude. We use a hybrid sigma-pressure coordinate system—near the surface, the levels are terrain following, while at higher altitudes they are defined as levels of constant pressure (Taylor et al., 2020). The EAMv2 model pairs the resolved atmospheric dynamical equations with a variety of the sub-grid parameterizations such as cumulus convection (G. J. Zhang & McFarlane, 1995), boundary layer cloud dynamics (Golaz et al., 2002), cloud micro-physics (Morrison & Gettelman, 2008), aerosol micro-physics and chemistry (Liu et al., 2016), and radiative transfer (Mlawer et al., 1997). The coarse-scaled simulations are run on an unstructured spherical element grid of approximately 1° (~ 110 [km]) resolution per sigma-level and 72 levels along the vertical direction, from 64[km], corresponding to ~ 0.1 [hPa] (level 1) down to the earth's surface (level 72). The vertical grid spacing is uneven, with the layer height ranging from 20 to 100 m near the surface up to 600 m in the upper atmosphere. We enforce appropriate boundary conditions over the Earth's surface in accordance with version 4.5 of the community land model (Oleson

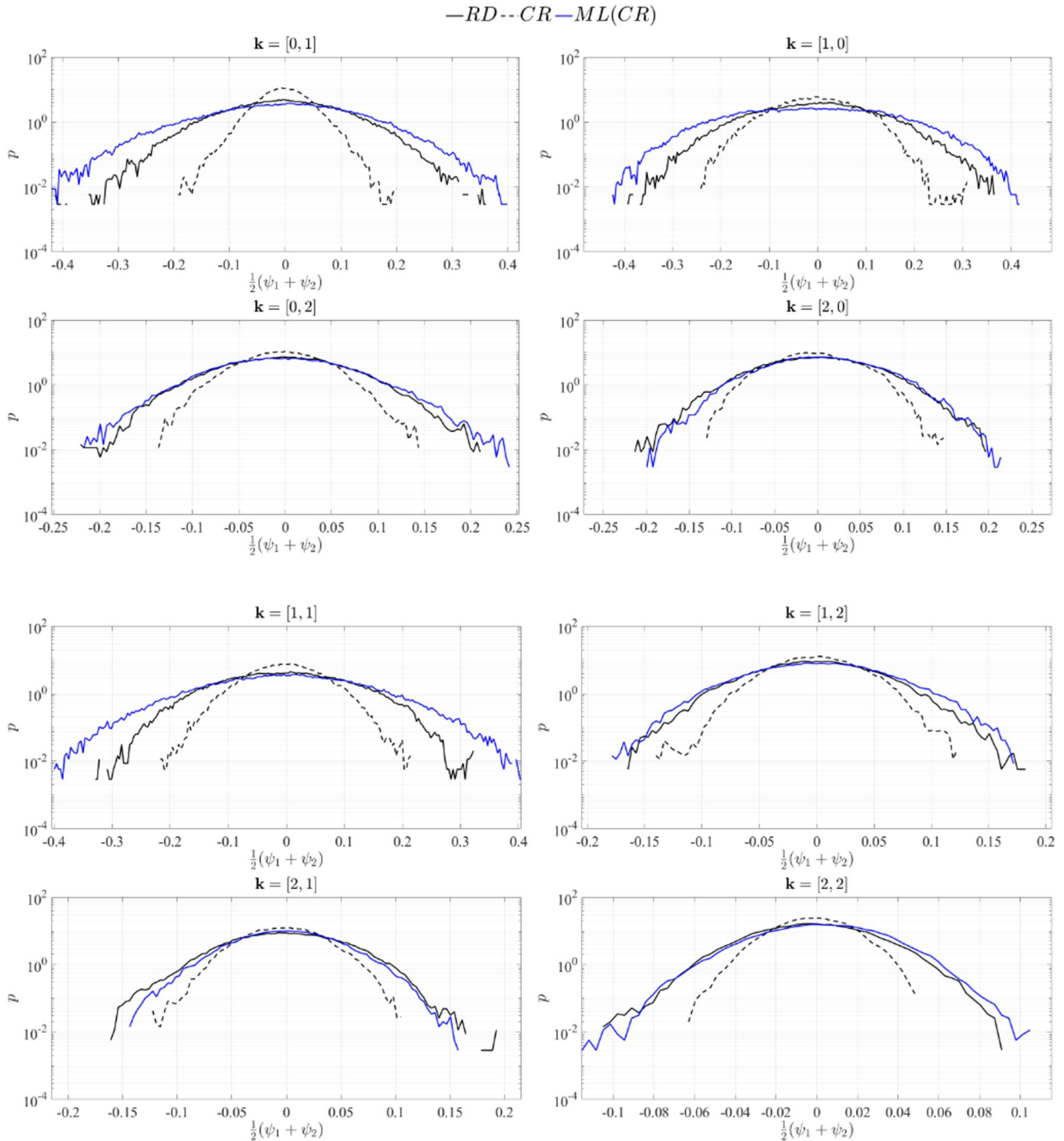


Figure 8. Probability density function of individual Fourier modes for $\beta = 1.1$ and $r = 0.5$. RD (black), ML(CR) (blue). Training data: $\beta = 2.0$ and $r = 0.1$.

et al., 2013). The (SST) and sea ice concentration boundary conditions are set according to the input4mip data sets (Reynolds et al., 2002).

In this case, the reference data used to generate the nudged training data and the validation reference is not a fully-resolved simulation but instead ERA5 reanalysis data (Hersbach et al., 2020) projected onto the coarse unstructured grid of EAMv2. The ERA5 data set combines observations with physics models to provide high-quality

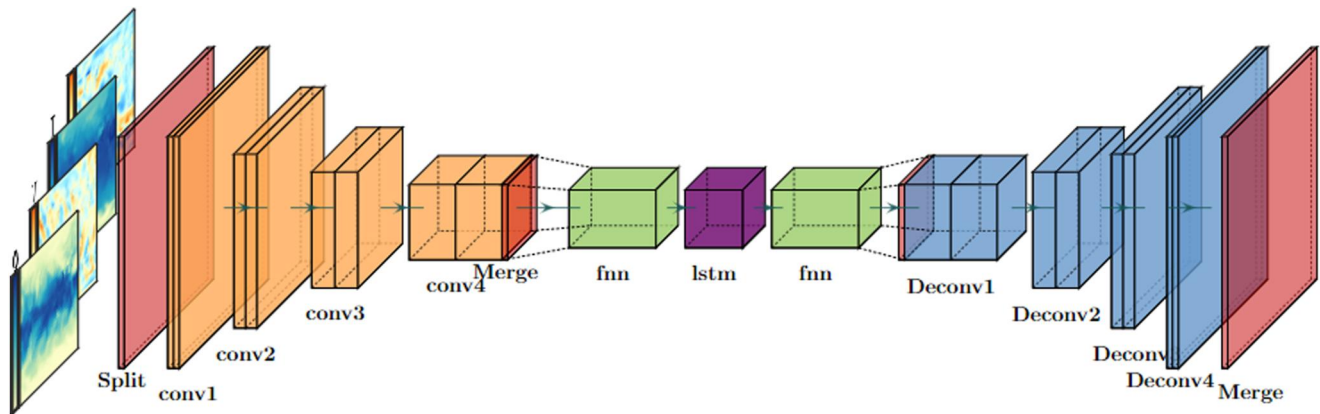


Figure 9. LSTM based neural network architecture used for the E3SM climate model.

reanalysis data on an hourly basis with a spatial resolution of 0.25° (~ 31 [km]). An outline of the practical implementation of the nudging is summarized in Appendix A1.

We do not perform any E3SM simulations at this fine resolution due to the prohibitive computational cost, and so in the following discussion any reference to E3SM data should be understood to represent the coarse model. Moving forward, the free-running data set will again be labeled as CR, the ML correction thereof as ML(CR), and the ERA5 reference data as RD. The data sets discussed herein contain information from 1979 to 2014, over which the climate system can be assumed to be in an approximately statistical steady state.

4.2. NN Architecture and Training Strategy

For the E3SM model we have developed a custom convolutional-LSTM hybrid network architecture. The architecture acts on a single sigma level, such that training is conducted for each level sequentially. The network receives as its input snapshots of the predictive variables $\mathbf{X} = \mathbf{X}(\theta, \phi, t, k)$ for fixed sigma level k . Afterward, a custom “split” layer separates the input into 25 non-overlapping subregions. These subregions are periodically padded via a custom padding process, tasked with respecting the spherical periodicity of the domain. Then, each subregion is independently passed through a series of four convolutional layers. The purpose of this process is to extract anisotropic local features in each subregion such as vapor transport.

Afterward, the local information extracted from each subregion is concatenated in a single vector via a custom “merge” layer. The global information is now passed through a linear fully-connected layer, that acts as a basis projection of the spatial data onto a reduced-order 20-dimensional latent space. The latent space data are then corrected by a LSTM layer (Hochreiter & Schmidhuber, 1997). Subsequently they are projected back to physical space via another linear fully-connected layer. Next, global information is split into the same subregions of the input, and distributed to another series of four independent deconvolution layers that upscale the data to the original resolution. Finally, a custom “merge” layer gathers the information from each subregion and produces the final corrected snapshot. A schematic of the configuration for training on a particular layer is shown in Figure 9.

The motivation behind using LSTM neural networks lies in their ability to incorporate (non-Markovian) memory effects into the reduced-order model. This ability stems from Takens embedding theorem (Takens, 1981). This theorem states that given delayed embeddings of a limited number of state variables, one can still obtain the attractor of the full system for the observed variables. In addition to temporal nonlocality, the model is nonlocal in space. Note, that in terms of the LSTM layer, this information comes in the form of the latent space coefficients, which in general correspond to global modes that correspond to rows of the fully connected layer’s matrix. Under the assumption that both fully-connected layers have linear activation functions, the model can be mathematically depicted as a basis projection. Hence, the fully connected layers act as projection schemes to (a) compress input data to a latent space of low dimensionality and (b) project the LSTM prediction to physical space. Such LSTM based models have been shown to be capable of improving predictions of reduced-order models in a variety of settings (Charalampopoulos & Sapsis, 2022; Harlim et al., 2021; Vlachas et al., 2018; Wan et al., 2018). However,

we note that other network architectures are possible, such as the recently proposed Fourier-Neural operators (Bonev et al., 2023; Guibas et al., 2022; Li et al., 2021, 2022) which have shown remarkable skill in data-driven weather prediction (Pathak et al., 2022).

The network is trained using a standard mean square error loss function

$$\mathcal{L} = \alpha \sum_t \sum_\phi \sum_\theta \cos\left(2\pi \frac{\theta}{360}\right) \|\mathbf{X}^{\text{ml}} - \mathbf{X}^{\text{rd}}\|^2, \quad (23)$$

where α is a normalization coefficient. As previously, training is performed using the nudged data set as input to the ML transformation. Each term in the sum is multiplied by a cosine that is a function of the latitude to showcase that the integration takes place over a sphere. If that term is absent, the model would over-emphasize on learning the corrections at the poles. Training was conducted over 1,000 epochs using data from the years 2007–2011, with the year 2012 used for validation during training.

4.3. Results

We apply our model to an unseen free-running coarse-scale simulations of the E3SM model (CR) over a 36 years horizon. These results are denoted as $ML(CR)$. The reference statistics used to evaluate our model predictions are computed from ERA5 reanalysis data over the years 1979–2014 and are denoted as RD . We also show the predictions of a free running E3SM simulation denoted CR , this serves as the baseline which our model is seeking to improve.

4.3.1. Global Statistics

First, we analyze the global 36-year statistics as a function of altitude, that is, for all sigma levels. In Figure 10, we show the time- and zonally-averaged biases for sigma-levels 10–72 of the simulations for (a–c) U , (e–g) T , (i–k) Q . We omit the highest sigma levels 1–10, as here the reference data is less reliable and thus obscures the analysis. The left column shows the biases of the free-running E3SM while the right column shows those of the ML corrected. The biases are normalized with the standard deviation of the quantity of interest for each sigma-level individually (sub-figures c, f, and i). For the case of Q for sigma-levels below $z = 35$, the standard deviation of level 35 was used for normalization. This is due to the fact that the values of Q in the upper atmosphere are extremely low and normalizing such errors by the standard deviation of their own sigma-level yielded very high biases for both predictions, making the metric misleading. The dotted regions indicate where the biases are statistically significant up to a 95% confidence level as quantified by a Student- t test. The ML correction notably corrects the strong overestimation of the specific humidity (bottom row) for sigma levels $z > 40$. The biases in temperature (middle row) in the upper atmosphere are also notably improved, however the improvement is less pronounced. In the case of the wind speed (top row), the ML correction does reduce the bias throughout the atmosphere, however, both the free running E3SM and the ML correction thereof retain significant biases in the upper atmosphere.

We now focus on the sigma level nearest the surface—level 72. Figure 11 shows the annual mean ERA5 reference data, as well as the biases of the free-running and ML corrected predictions. The ML correction reduces the global RMSE by 18%, 19%, and 36% for U , T , and Q respectively. Regionally, the benefits of our model correction are best seen in the equatorial and south polar regions. In the former, the free-running solution significantly overestimates the specific humidity, while the ML correction is relatively free of any such systematic bias. Then in the latter, the uncorrected simulation significantly underestimates the temperature, a deficit which is remedied with the ML correction. To illustrate the temporal evolution of the near surface biases we also show in Figure 12 the time versus latitude Hovmöller diagrams of the monthly mean zonal mean bias in U , T , and Q over the time period 1979–2014. We note that the period 2007–2014 is part of our training data. Consistent with the results in Figure 11, our ML correction consistently reduces the zonal mean biases of all three quantities. The most significant improvements are observed in T and Q , for which the performance of the ML correction is greatest in the tropical and subtropical regions. Furthermore, in those regions where we observe significant bias reduction, the corrections persist robustly across the years outside the training period. However, there is an over-correction of the positive biases in Q in the tropical regions during the period 1979–2002 (Figure 12c). This is possibly because

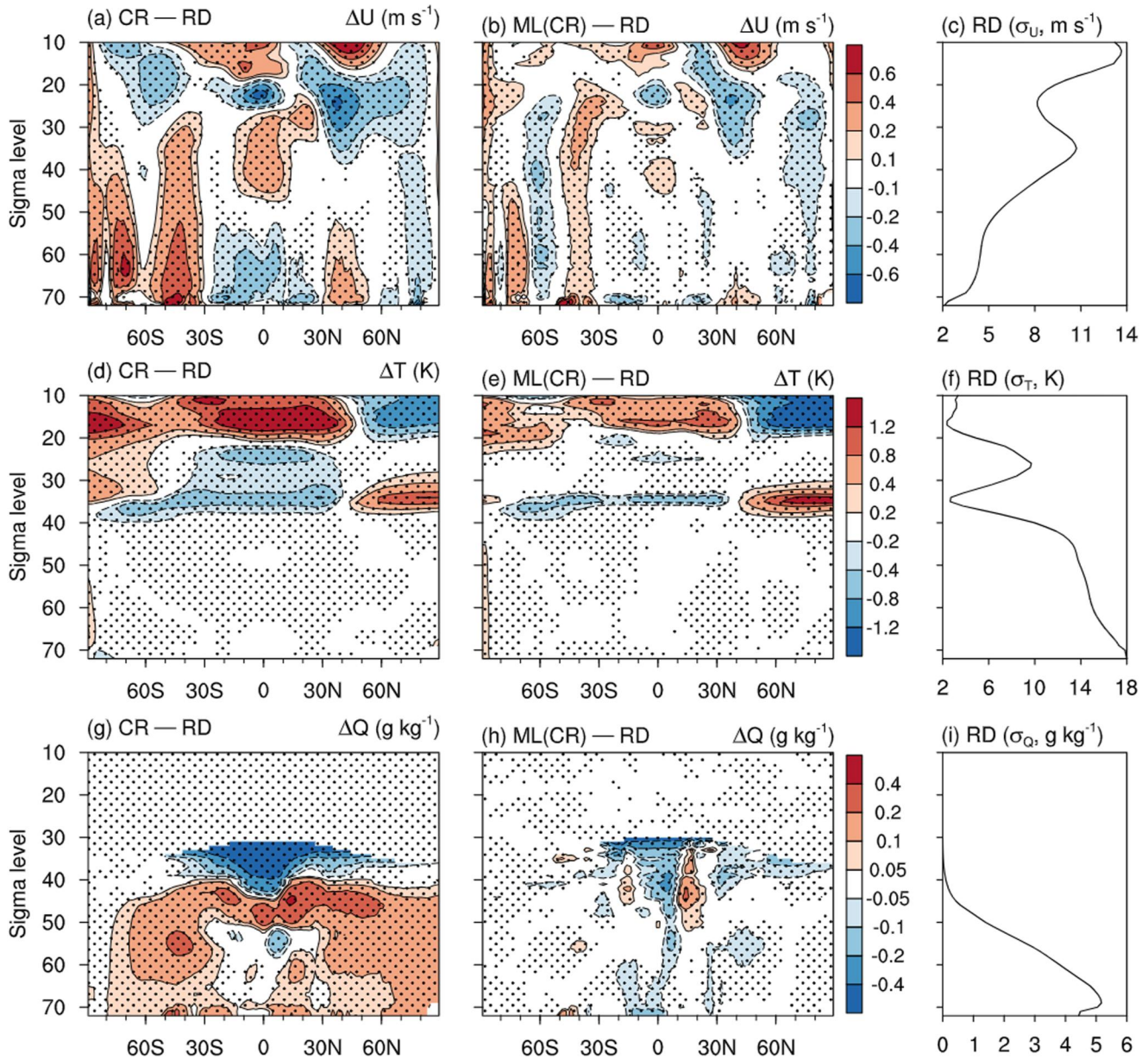


Figure 10. Zonally-averaged 36-year annual mean biases for all sigma-level of the simulations, for normalized zonal velocity U (a–c), temperature T (e–g), and specific humidity Q (i–k). Free running coarse E3SM simulation (CR) (left) and ML-correction (ML(CR)) (right). Standard deviation σ of each quantity at the specific sigma-level shown (d, h, and i).

the training data is too short to capture the multi-decade trend of the E3SM model increasingly overestimating the humidity in the tropics.

Figure 13 shows the aggregate probability density function at sigma level 72 across the globe for the same 36 year period. The probability density functions are computed using the 36×12 monthly mean values at each grid point. The ML correction significantly improves the predicted distributions in wind speed U , V (a, b) and specific humidity Q (d). Critically, the improvements are most pronounced in the tails of the distribution, which are critical for quantifying the risks of extreme weather events. There is very little improvement in the temperature (T), however, in this case the E3SM prediction alone is already quite accurate.

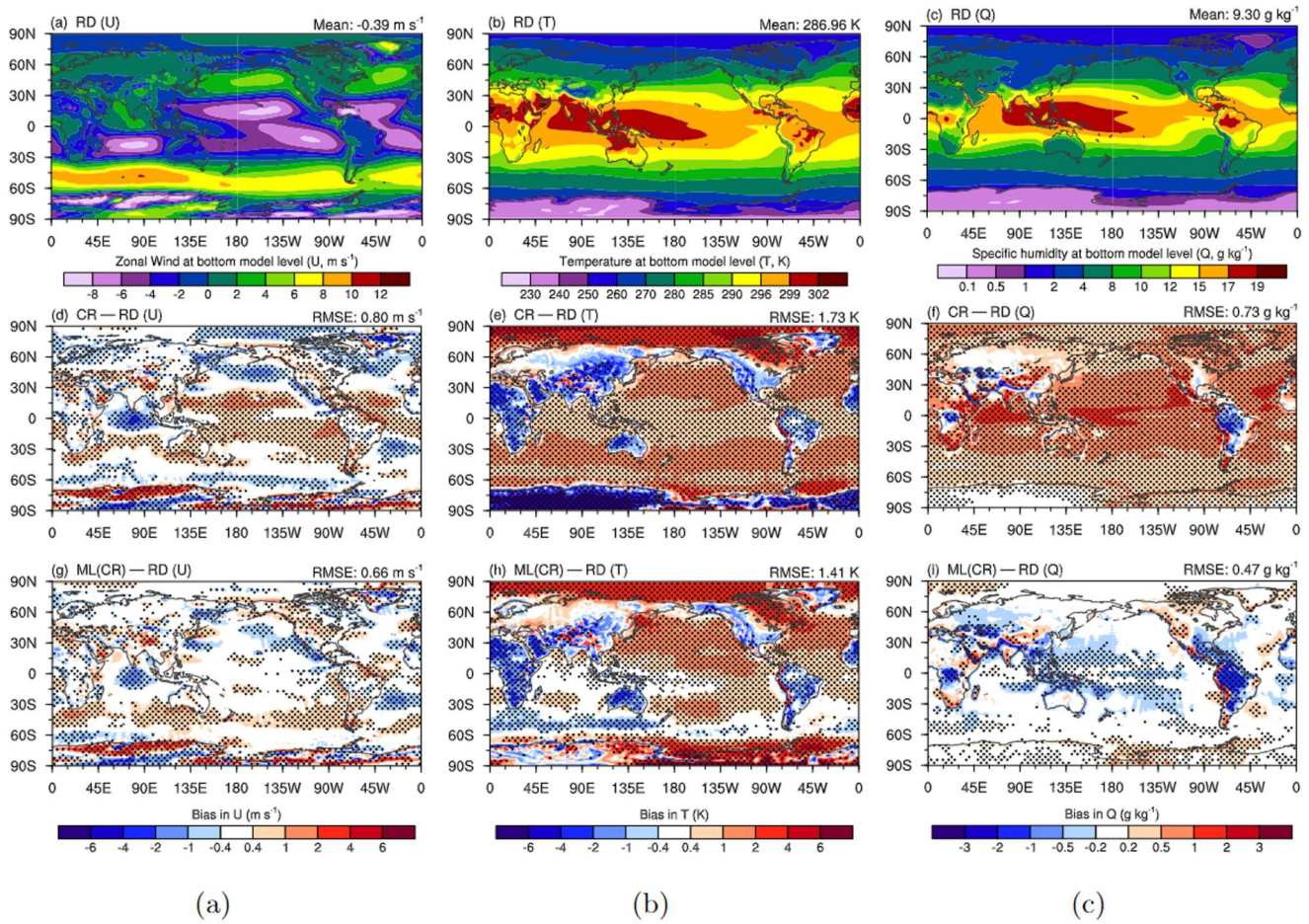


Figure 11. Global 36-year time-averaged biases at the lower-most sigma-level with respect to ERA5 for time-averaged zonal velocity U , temperature T and specific humidity Q . Top row corresponds to the reference data (RD), second row corresponds to a free-running simulation (CR) and bottom row corresponds to ML-correction (ML(CR)).

4.3.2. IVT

We now move to predict statistics for a derived integral quantity, the mean IVT. The IVT quantifies the vertically integrated mass transport of water vapor and is defined as

$$IVT(t, \theta, \phi) \equiv \sqrt{IVT_U^2 + IVT_V^2} \quad (24)$$

where IVT_U and IVT_V are the east-west and north-south components defined as

$$IVT_U(t, \theta, \phi) \equiv \frac{1}{g} \int Q(t, \theta, \phi, p) U(t, \theta, \phi, p) dp \quad (25)$$

and similarly for IVT_V , and where the vertical coordinate has been re-parameterized in terms of pressure. Regions of concentrated IVT are known as atmospheric rivers (AR) and are associated with heavy precipitation and a variety of extreme weather events—both beneficial and detrimental. For example, on the open ocean, ARs are generally associated with extratropical cyclones, and upon landfall ARs have the potential to alleviate drought conditions or lead to significant storm damage (Payne et al., 2020). Therefore, the ability to correctly predict the statistics of the IVT—and thus ARs—is a crucial metric by which to evaluate our ML correction operator. Although it is beyond the scope of this work, the interested reader is referred to S. Zhang et al. (2023) for a

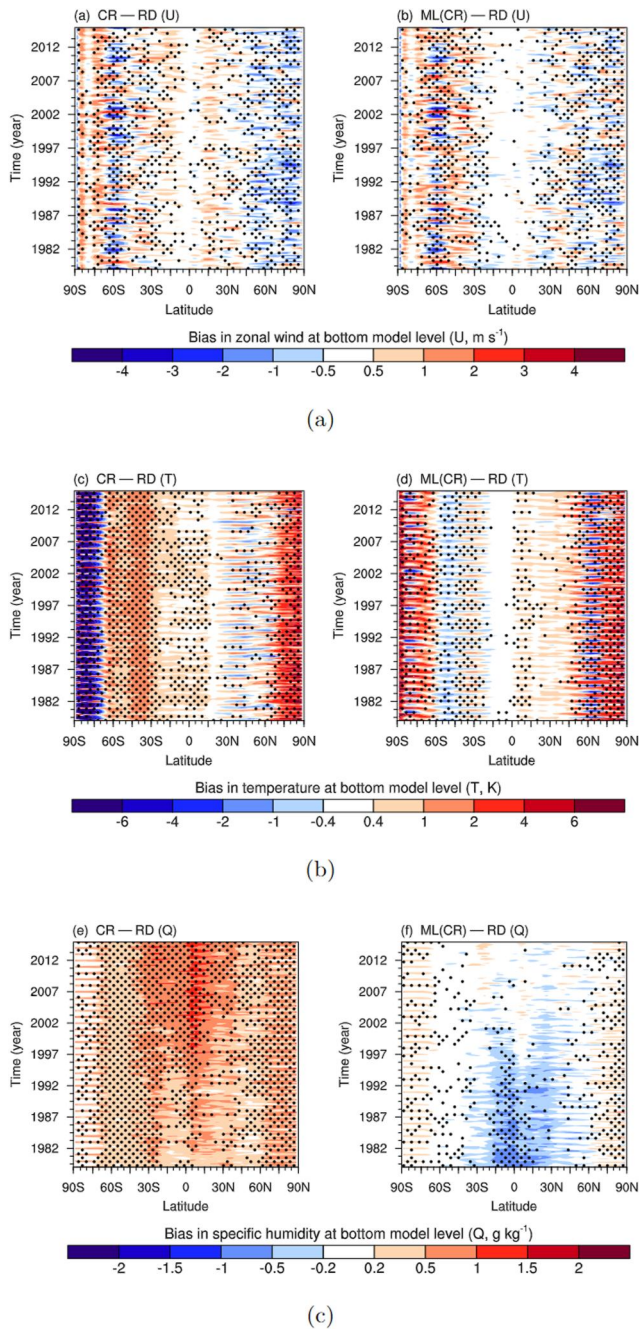


Figure 12. Hovmöller Diagrams of biases at the lower-most sigma-level with respect to ERA5 for time-averaged zonal (a) velocity U , (b) temperature T , and (c) specific humidity Q . Free running coarse E3SM simulation (CR) (left) and ML-correction (ML(CR)) (right).

5. Discussion

We have introduced a method to machine learn correction operators to improve the statistics of under-resolved simulations of turbulent dynamical systems. The premise of the proposed strategy is to generate training data pairs which are minimally affected by chaotic divergence. Instead of using an *arbitrary* coarse trajectory as the training

detailed discussion of our method applied to the statistics of other extreme climate events such as tropical cyclones.

From a ML point of view, accurately predicting the spatial features of extreme events, which are quantified by highly anisotropic quantities such as IVT, requires accurately mapping local flow features between the under- and fully-resolved trajectories. It is for this reason, that we have implemented the domain-splitting and local convolution layers in the network architecture described in Section 4.2.

In Figure 14, we show the 36-year annual mean of the IVT across the globe. The top figure corresponds to the ERA5 reanalysis data, and below that are the biases of the free-running E3SM simulation, as well as the machine learned correction. Overall, the ML correction decreases the global root-mean-square error (RMSE) by 51% compared to the free-running E3SM solution. Furthermore, the ML correction significantly decreases several systematic regional biases throughout the domain. Note for example, that the ML significantly reduces the strong positive bias of the free-running E3SM simulation over Southeast Asia and in the southern oceans around 45° of latitude.

4.3.3. Regional Statistics

In addition to global statistics, policy makers preparing for the increased risks of climate change require accurate risk analysis over a range of spatial scales. Therefore we also analyze the statistics of the predicted climate over several regions of varying size: the tropics, mid-latitude, continental US, northeast US, northern Europe, and the northwest Pacific. The size and location used in the following results are summarized in Table 1. As in Section 4.3.1 we focus on sigma level 72, the level closest to the surface. Figures 15–17 show the probability density functions of the four progress variables U , V , T , and Q in the tropics, mid-latitude, and the northwest Pacific regions. Result for the remaining regions are included in Appendix A2. The reanalysis reference is shown in solid black, the free-running E3SM and ML correction thereof are shown in dashed black and blue respectively. Again, we see that the ML correction is most pronounced in regions where the E3SM model alone is most biased. Most notably the specific humidity Q (subplot d in Figures 15–17) and meridional wind speed (V) (subplot b in Figures 15–17) where for all regions the ML correction brings the tails of the predicted distribution into good agreement with ERA5 data. See also Figure 15a, where the ML correction does significantly improve the prediction of the zonal wind speed (U). As with the global statistics, the ML correction has only minor impacts on the distributions of temperature (T). However, with the exception of the tropics region (Figure 15c) this is generally well predicted by the E3SM model alone and notably in no region does our ML correction significantly increase bias. The fact that our correction operator is able to improve predictions across all variables and over a range of spatial scales is a promising result, as it shows that the predicted flow field could in principle be further used for targeted super-resolution to predict local features on scales smaller than the grid of the coarse model.

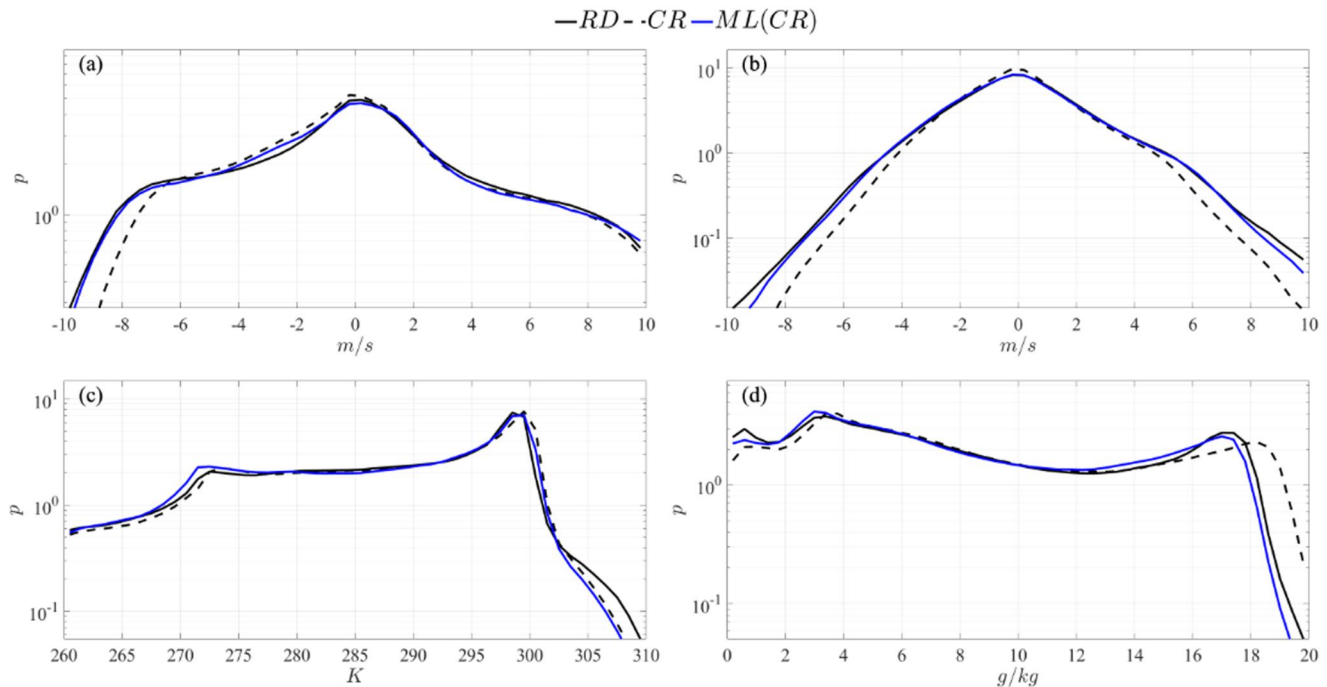


Figure 13. Global 36-year probability density function for surface sigma-level 72. (a) U , (b) V , (c) T , (d) Q . Results are shown for ERA5 reanalysis data (RD) (solid black), free-running data (CR) (dashed black), and machine learning corrections (ML(CR)) (blue).

input, we used a coarse trajectory *nudged* toward the training target trajectory. This nudged trajectory predominately obeys the dynamics of the coarse model, yet is constrained from *randomly wandering* too far from the reference. In essence, it is an approximation of the one (of infinitely many) trajectory of the coarse model which is closest to the reference data. Once trained on this specific pair of trajectories, an ML operator can reliably map *any* free-running coarse trajectory into the attractor of the reference data. The critical benefit of such an operator is that it acts on data in a post-processing manner, and is thus unaffected by the stability issues, and practical implementation challenges, which plague machine learned corrections of the system dynamics.

A key aspect of the proposed approach is the ability to incorporate, directly into the learning process, dynamical information that goes beyond statistics of the training data. This is achieved through an objective function that is matching *trajectories* rather than their statistics. This is critical especially for extreme events, where the key information “lives” in the very structure of the trajectory over the short duration of such events. Cost functions formulated to match statistics, either need to incorporate high order statistical information (something that is practically impossible because of both inadequate data but also vast computational cost) or they are doomed to have poor generalization properties since low order statistics (e.g., spectrum) cannot “see” the dynamics of extreme events. On the other hand, the formulated approach eliminates the divergence due to chaotic behavior and uses the maximum information from the reference data by training *in the time domain*, that is, directly fixing the structure of the trajectory near an extreme event. This allows for unprecedented improvement especially for extreme event statistics.

The proposed strategy was first illustrated on a prototypical two layer quasi-geostrophic climate model using a simple LSTM network architecture. In this reduced order system our ML correction operator was able to bring the global, and scale-by-scale statistics of a severely under-resolved simulation, simulated on a 24×24 grid, into good agreement with the fully-resolved reference solved on a 128×128 grid. Additionally, we demonstrated the ability to accurately predict statistics for time horizons much longer than the training data, and for parameter regimes outside of that training data. We then applied our framework to a realistic climate model—the Energy Exascale Earth System Model (E3SM) solved on a grid with approximately 110 km horizontal resolution. In this case, the reference data used as the training target and the evaluation metric was not a fully resolved simulation, but ERA5 reanalysis data. To address this far more complex system, we designed a network architecture which combined the LSTM base we used for the simpler QG system with overlapping convolutional layers used to

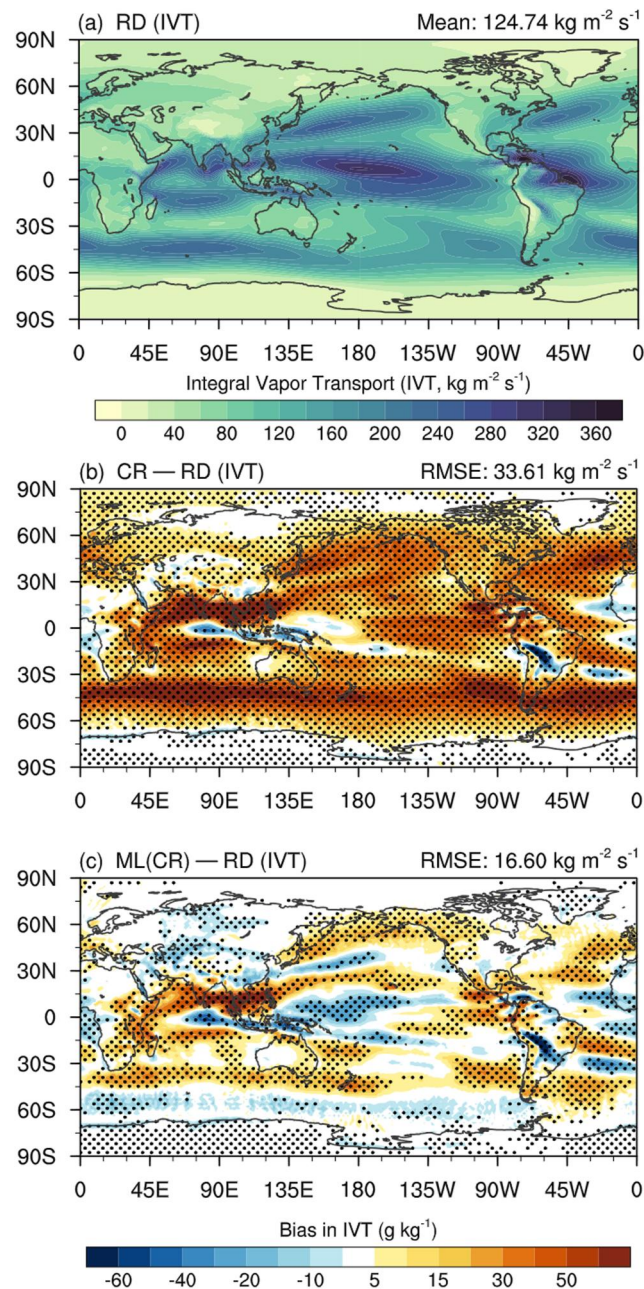


Figure 14. 36-year annual mean integrated vapor transport predictions. From top to bottom, ERA5 reference, free-running E3SM bias, machine learning correction bias.

Table 1
Summary of Regions Analyzed in Section 4.3.3

Region	Latitude	Longitude
Mid-latitude	30°–60°S and 30°–60°N	0°–360°
Tropics	20°S–20°N	0°–360°
Continental US	25°–55°N	90°–120°W
Northeastern US	25°–55°N	60°–90°W
Northern Europe	40°–70°N	10°–40°E
Northwest Pacific	30°–60°N	150°–180°E

extract local anisotropic features from the input data. We found that our ML correction significantly reduced the bias of the E3SM solution, bringing the statistics of the wind speeds and specific humidity into good agreement with reanalysis data on both a global and regional level. The debiasing capabilities of our ML correction were less pronounced in the case of temperature, for which the improvements, especially in the tails of the distributions were more modest, and more region dependent. The improvement in the wind speed and humidity statistics however are especially notable as these variables were not well approximated by the free-running E3SM solution. In particular, the correction operator significantly improved the predictions of the tails of these

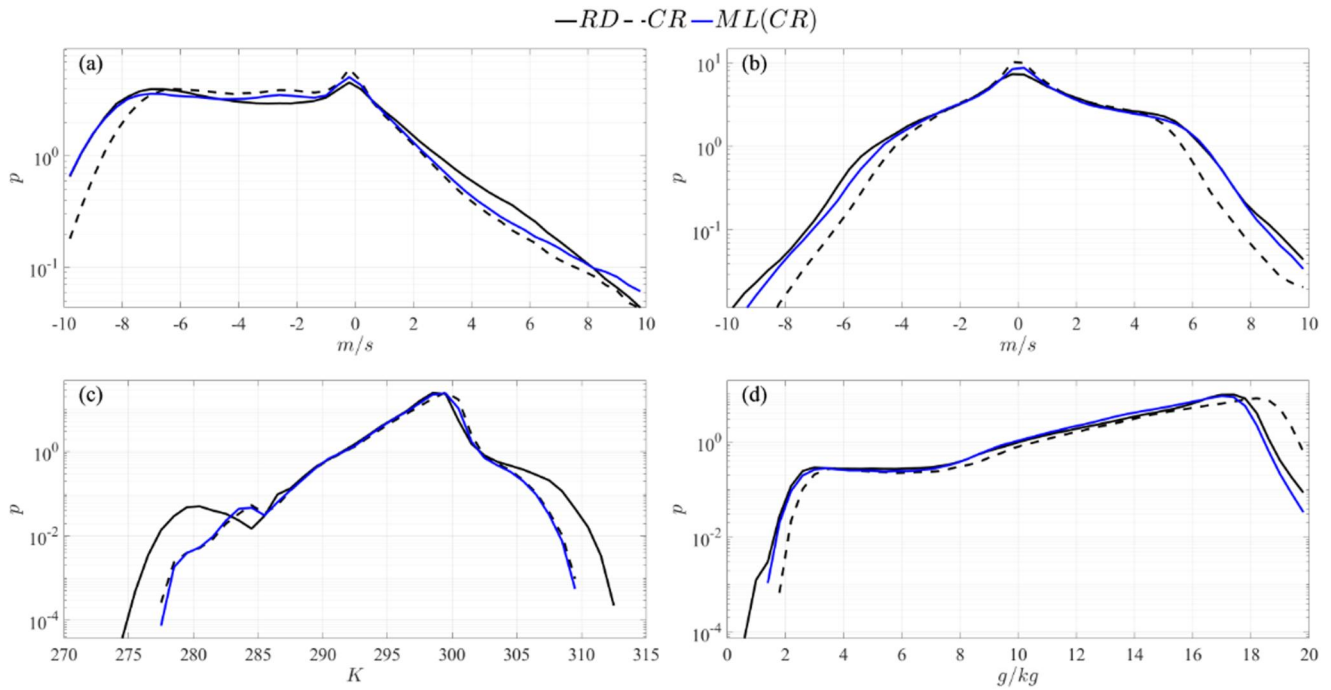


Figure 15. 36-year probability density function for surface sigma-level 72 in the tropics. (a) U , (b) V , (c) T , (d) Q . Results are shown for ERA5 reanalysis data (RD), free-running data (CR), and machine learning corrections.

distributions which are critical for quantifying the risks of extreme weather events. In addition to the primitive variables, we also analyzed the mean IVT, a highly anisotropic integral quantity of particular practical interest as it drives AR and thus precipitation. Here the improved predictions in the wind speed and humidity of our ML correction combined to reduce the overall RMSE in IVT by 51%, and successfully removed several systematic

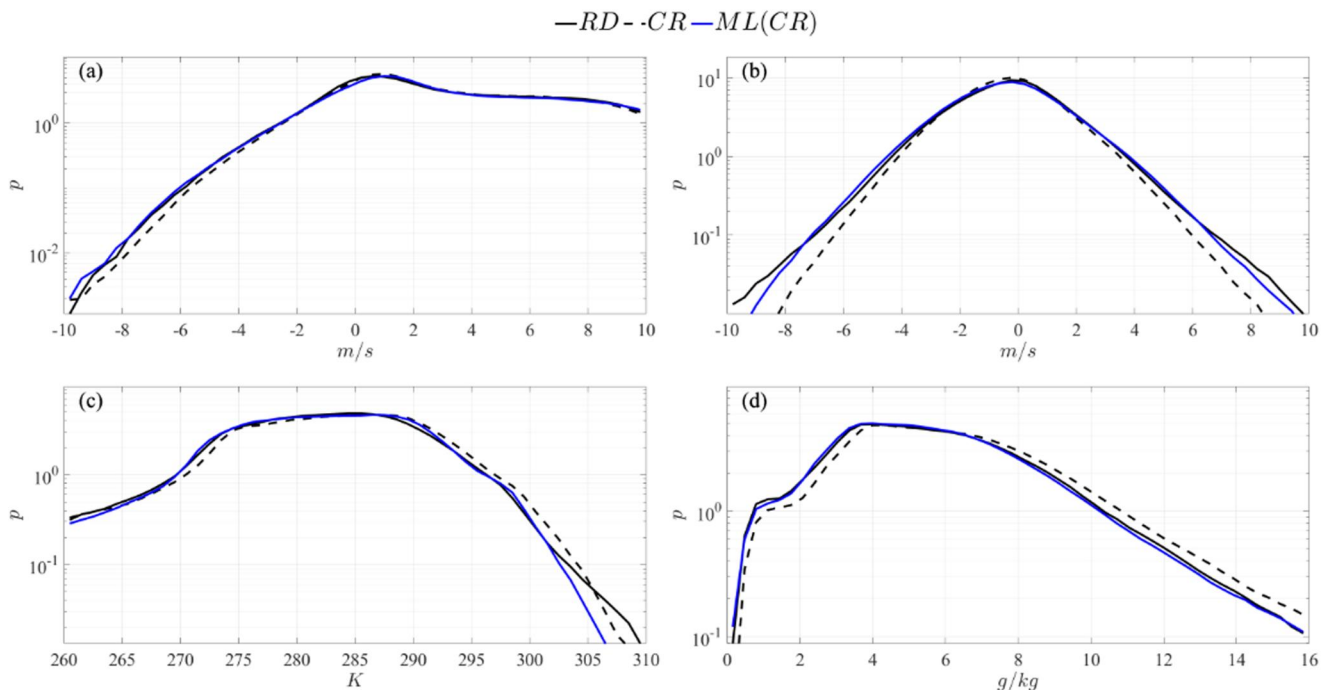


Figure 16. 36-year probability density function for surface sigma-level 72 in the mid-latitude region. (a) U , (b) V , (c) T , (d) Q . Results are shown for ERA5 reanalysis data (RD) (solid black), free-running data (CR) (dashed black), and machine learning corrections (blue).

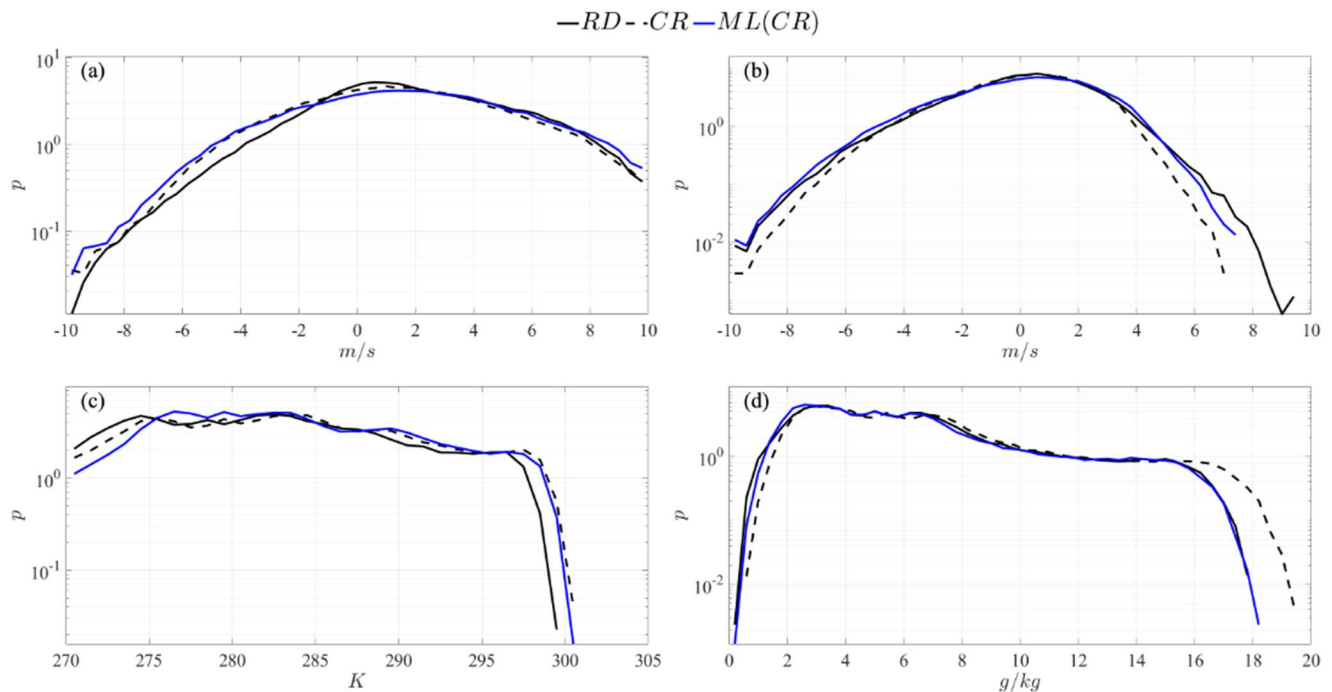


Figure 17. 36-year probability density function for surface sigma-level 72 in northwest Pacific. (a) U , (b) V , (c) T , (d) Q . Results are shown for ERA5 reanalysis data (RD) (solid black), free-running data (CR) (dashed black), and machine learning corrections (blue).

regional biases of the coarse model, such as its tendency to underpredict the vapor transport in the southern hemisphere.

While the proposed methodology was demonstrated to be effective for the prediction of a multitude of climate metrics, some limitations of the current setup should be stated. First, the approach works well under the assumption that the climate is in a statistically steady state, for which a mapping can be learned through the proposed training scheme. Hence, applying the learned model in situations where the climate undergoes a transitory phase may hinder its performance, unless similar transitory intervals are included in the training data. This is particularly true if the transition is not captured at all by the coarse-scale model. Furthermore, when applied to future climate scenarios with drastically different forcing, the requirement for reference data—which may not be available at high resolution for long times—makes it difficult to assess the predictive powers of our approach a priori. For such runs to be included in training, high-fidelity simulations would have to be used as reference and the coarse models nudged toward them. This limitation however is true for online data-driven correction schemes as well since most such models lack concrete error bounds for out-of-sample predictions. Furthermore, for the application of the scheme to dynamical systems broadly, there is no guarantee that a nudged simulation exists that follows the reference data closely while satisfying the dynamics of the coarse simulation. Essentially, if the coarse model is too far from the reference data, that is, too under-resolved or neglecting too much important physics there is no guarantee the process will work.

One of the main advantages of the proposed framework is its generality and non-intrusive nature. Theoretically, intrusive online approaches act on the dynamics of the system, but practically, this means they act on *software*, that is, they must be integrated with existing code stacks. For modern ESMs, this code stack can be complex or proprietary, making the implementation of such strategies difficult or even impossible if the source code is unavailable. On the other hand, non-intrusive approaches, such as the one proposed here, act on *data*—meaning the model is agnostic to the specific software implementation of the model generating the data. Generating the training data does require implementing a nudging tendency in the climate model code, however, this is generally a much less invasive task than integrating an ML operator, which may be implemented in a different software language than the climate model itself (J. McGibbon et al., 2021). Then once trained the model can be used without further intrusion into the core ESM. Another strength, is that the proposed framework provides predictions of all progress variables, (U , V , T , Q), at all grid points and all sigma levels—a feature not shared by

all debiasing schemes. This in turn means that the flow fields predicted by our correction operator could then be used for local super-resolution (down-scaling) to investigate local climate forecasting and impact assessment. However, further work is required to investigate the ability of our approach to improve the statistics of other climate metrics such as precipitation and to ensure that the corrected fields obey basic physical constraints such as geostrophic balance or conservation of mass and energy over the spatio-temporal scales relevant to such local analysis. We believe that by lowering these barriers to adoption, our approach has the potential to significantly accelerate and democratize the implementation of data-driven climate modeling. To this end, extensions of our approach such as built in uncertainty quantification, physics informed constraints, and grid-agnostic network architectures—which could allow for applications across different ESMs—are the topic of ongoing research.

Appendix A

A1. Nudging Implementation in E3SM

Here we briefly outline the practical implementation of the nudging strategy in the E3SM model used to train the ML correction operator used to generate the results in Section 4. We follow the formulation of Sun et al. (2019) and Zhang et al. (2022), for which the nudged governing equations of the E3SM model takes the form

$$\frac{\partial \mathbf{X}}{\partial t} = \underbrace{\mathbf{D}(\mathbf{X})}_{\text{dynamics}} + \underbrace{\mathbf{P}(\mathbf{X})}_{\text{physics}} - \underbrace{\mathcal{N}(\mathbf{X}, \mathbf{X}^{RD})}_{\text{nudging}} \quad (\text{A1})$$

where \mathbf{D} represents the resolved dynamics, \mathbf{P} represents the parameterized physics and \mathcal{N} is the nudging tendency. The nudging tendency is applied at each grid point and is specifically implemented as

$$\mathcal{N}(\mathbf{X}, \mathbf{X}^{RD}) = \begin{cases} 0, & \text{if } P \leq 1 \text{ Pa} \\ -\frac{\mathbf{X} - \mathbf{X}^{RD}}{\tau} \times \frac{P_m}{P_0}, & \text{if } 1 \text{ Pa} < P \leq P_0 \\ -\frac{\mathbf{X} - \mathbf{X}^{RD}}{\tau} \times \frac{1}{2} \left[1 + \tanh\left(\frac{Z - Z_b}{0.1Z_b}\right) \right], & \text{if } Z \leq Z_p \\ -\frac{\mathbf{X} - \mathbf{X}^{RD}}{\tau}, & \text{otherwise} \end{cases} \quad (\text{A2})$$

where $\mathbf{X} = (U, V, T, Q)$ is the state variable, \mathbf{X}^{RD} is the ERA5 reference, P_m and Z_m represent the atmospheric pressure and geopotential height at a given sigma level, and τ denotes the relaxation time scale. Following Sun et al. (2019) and Zhang et al. (2022) we fix $\tau = 6$ hr. The simulation uses a time step of 0.5 hr and the ERA5 reference data is defined at 3-hourly increments and interpolated at each time step using the linear temporal interpolation described in Sun et al. (2019). The quantities P_0 and Z_b are user defined threshold parameters which govern how the nudging tendency is modulated in the upper and lower ends of the atmosphere. Z_b is set at the planetary boundary layer height, which is diagnosed and dynamically set at each time step. P_0 is set to 30, 30, 10, and 100 Pa for the variables U, V, T, Q respectively and held constant throughout the simulation. This modulation in the upper and lower sigma levels differs from the default formulation proposed by Sun et al. (2019) and Zhang et al. (2022), however, it is implemented here to account for uncertainties in our specific reference data. We de-emphasize the nudging tendency in the upper atmosphere due to the deteriorating quality of the ERA5 reanalysis data at those altitudes, while near-surface the concern is the significant errors which arise over the high-terrain regions when ERA5 data is mapped onto the E3SM model grid.

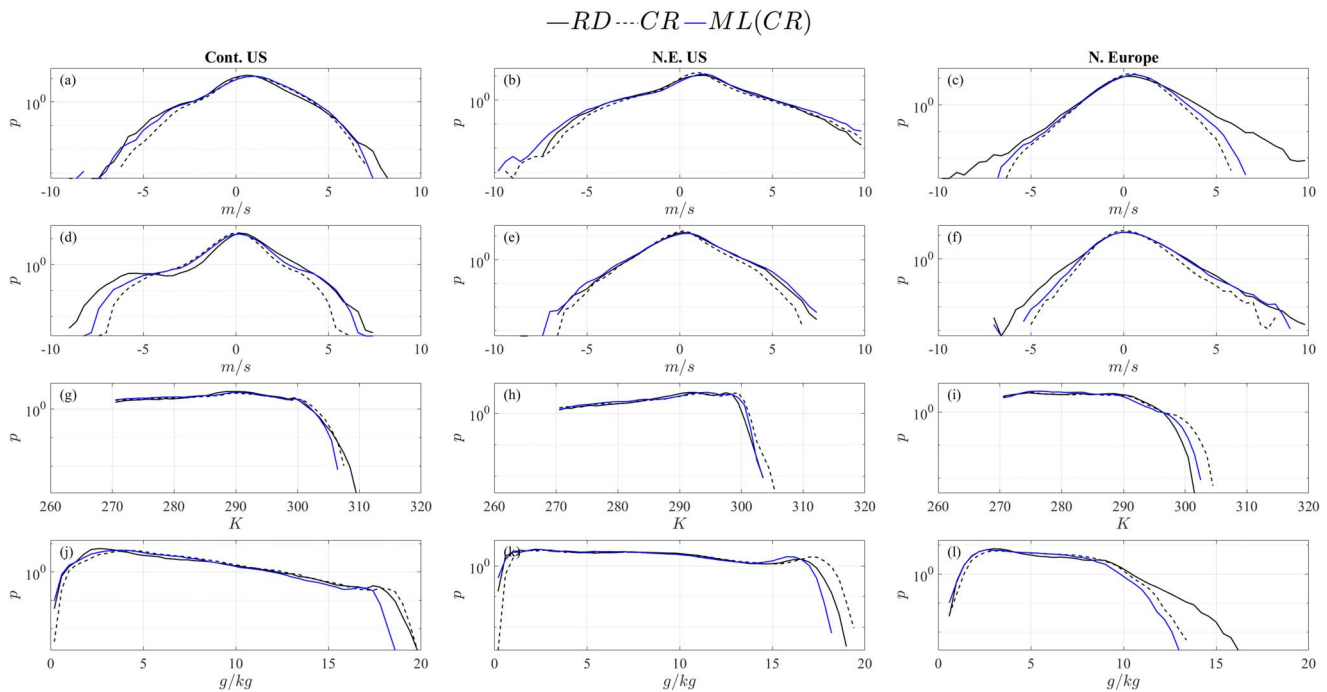


Figure A1. 36-year probability density function for surface sigma-level 72 for Continental US (left column), northeastern US (center column) and northern Europe (right column). (a–c) U , (d–f) V , (g–i) T , and (j–l) Q . Results are shown for ERA5 reanalysis data (RD) (solid black), free-running data (CR) (dashed black), and machine learning corrections (blue).

A2. Additional E3SM Results

Here we show some additional results for Section 4. Figure A1 shows the regional probability density functions for the regions not shown in Section 4: Continental US (left column), northeastern US (center column) and northern Europe (right column) at the surface sigma level 72.

Data Availability Statement

The source code for the E3SM (E3SM Project, 2021) climate model used to generate the simulations discussed in Section 4 was obtained from the Energy Exascale Earth System Model project, sponsored by the U.S. Department of Energy, Office of Science, Office of Biological and Environmental Research. The ERA5 reanalysis data used as a reference for training the ML model and generating the reference data in Section 4 is available at the Copernicus Climate Change Service (C3S) Climate Data Store via <https://doi.org/10.24381/cds.bd0915c6> (? , ?). The software and data needed to generate the results described here can be found on Zenodo at <https://zenodo.org/doi/10.5281/zenodo.10657047> (Barthel et al., 2023).

Acknowledgments

This research has been supported by the DARPA Grant HR00112290029 under the program “AI-Assisted Climate Tipping Point Modeling” supported by the Program Manager Dr. Joshua Elliott. Pacific Northwest National Laboratory is operated for the U.S. Department of Energy by Battelle Memorial Institute under Contract DE-AC05-76RL01830. Computational resources for the material shown in this work were provided by Anvil super computer through the ACCESS program. The authors thank Prof. G. Karniadakis for stimulating discussions on this work. We are also grateful to Dr. S. Khurshid for advise and support on using the Anvil super computer.

References

- Allen, S., Barros, V., Canada, I., UK, D., Cardona, O., Cutter, S., et al. (2012). Managing the risks of extreme events and disasters to advance climate change adaptation. Special Report of Working Groups I and II of the Intergovernmental Panel on Climate Change (Technical Report). <https://doi.org/10.13140/2.1.3117.9529>
- Arbabi, H., & Sapsis, T. (2022). Generative stochastic modeling of strongly nonlinear flows with non-Gaussian statistics. *SIAM/ASA Journal on Uncertainty Quantification*, *10*(2), 555–583. <https://doi.org/10.1137/20M1359833>
- Arcomano, T., Szunyogh, I., Wikner, A., Hunt, B. R., & Ott, E. (2023). A hybrid atmospheric model incorporating machine learning can capture dynamical processes not captured by its physics-based component. *Geophysical Research Letters*, *50*(8), e2022GL102649. <https://doi.org/10.1029/2022GL102649>
- Arcomano, T., Szunyogh, I., Wikner, A., Pathak, J., Hunt, B. R., & Ott, E. (2022). A hybrid approach to atmospheric modeling that combines machine learning with a physics-based numerical model. *Journal of Advances in Modeling Earth Systems*, *14*(3), e2021MS002712. <https://doi.org/10.1029/2021MS002712>
- Barriopedro, D., Fischer, E. M., Luterbacher, J., Trigo, R. M., & García-Herrera, R. (2011). The hot summer of 2010: Redrawing the temperature record map of Europe. *Science*, *332*(6026), 220–224. <https://doi.org/10.1126/science.1201224>
- Barthel, B., Zhang, S., Charalampopoulos, A.-T., & Themistoklis, S. (2023). Analysis scripts and dataset for Barthel et al. (2023) [Dataset]. Zenodo. <https://doi.org/10.5281/zenodo.10120306>

- Bauer, P., Stevens, B., & Hazeleger, W. (2021). A digital twin of Earth for the green transition. *Nature Climate Change*, *11*(2), 80–83. <https://doi.org/10.1038/s41558-021-00986-y>
- Bevacqua, E., Suarez-Gutierrez, L., Jézéquel, A., Lehner, F., Vrac, M., Yiou, P., & Zscheischler, J. (2023). Advancing research on compound weather and climate events via large ensemble model simulations. *Nature Communications*, *14*(1), 2145. <https://doi.org/10.1038/s41467-023-37847-5>
- Blanchard, A., Parashar, N., Dodov, B., Lessig, C., & Sapsis, T. (2022). A Multi-scale deep learning framework for projecting weather extremes. In *Climate Change AI*. Retrieved from <https://www.climatechange.ai/papers/neurips2022/65>
- Bloom, A. A., Exbrayat, J.-F., van der Velde, I. R., Feng, L., & Williams, M. (2016). The decadal state of the terrestrial carbon cycle: Global retrievals of terrestrial carbon allocation, pools, and residence times. *Proceedings of the National Academy of Sciences*, *113*(5), 1285–1290. <https://doi.org/10.1073/pnas.1515160113>
- Bonev, B., Kurth, T., Hundt, C., Pathak, J., Baust, M., Kashinath, K., & Anandkumar, A. (2023). Spherical Fourier neural operators: Learning stable dynamics on the sphere. Retrieved from <https://arxiv.org/abs/2306.03838v1>
- Bora, A., Shukla, K., Zhang, S., Harrop, B., Leung, R., & Karniadakis, G. E. (2023). Learning bias corrections for climate models using deep neural operators. arXiv. <https://doi.org/10.48550/arXiv.2302.03173>
- Brenowitz, N. D., & Bretherton, C. S. (2019). Spatially extended tests of a neural network parametrization trained by coarse-graining. *Journal of Advances in Modeling Earth Systems*, *11*(8), 2728–2744. <https://doi.org/10.1029/2019MS001711>
- Bretherton, C. S., Henn, B., Kwa, A., Brenowitz, N. D., Watt-Meyer, O., McGibbon, J., et al. (2022). Correcting coarse-grid weather and climate models by machine learning from global storm-resolving simulations. *Journal of Advances in Modeling Earth Systems*, *14*(2), e2021MS002794. <https://doi.org/10.1029/2021MS002794>
- Buchta, D. A., & Zaki, T. A. (2021). Observation-infused simulations of high-speed boundary-layer transition. *Journal of Fluid Mechanics*, *916*, A44. <https://doi.org/10.1017/jfm.2021.172>
- Charalampopoulos, A.-T. G., & Sapsis, T. P. (2022). Machine-learning energy-preserving nonlocal closures for turbulent fluid flows and inertial tracers. *Physical Review Fluids*, *7*(2), 024305. <https://doi.org/10.1103/PhysRevFluids.7.024305>
- Clark, S. K., Brenowitz, N. D., Henn, B., Kwa, A., McGibbon, J., Perkins, W. A., et al. (2022). Correcting a 200 km resolution climate model in multiple climates by machine learning from 25 km resolution simulations. *Journal of Advances in Modeling Earth Systems*, *14*(9), e2022MS003219. <https://doi.org/10.1029/2022MS003219>
- Dennis, J. M., Edwards, J., Evans, K. J., Guba, O., Lauritzen, P. H., Mirin, A. A., et al. (2012). CAM-SE: A scalable spectral element dynamical core for The Community Atmosphere Model. *The International Journal of High Performance Computing Applications*, *26*(1), 74–89. <https://doi.org/10.1177/1094342011428142>
- E3SM Project, DOE. (2021). Energy exascale earth system model v2.0 [Computer Software]. *U.S. Department of Energy*. <https://doi.org/10.11578/E3SM/dc.20210927.1>
- Fiedler, T., Pitman, A. J., Mackenzie, K., Wood, N., Jakob, C., & Perkins-Kirkpatrick, S. E. (2021). Business risk and the emergence of climate analytics. *Nature Climate Change*, *11*(2), 87–94. <https://doi.org/10.1038/s41558-020-00984-6>
- Fischer, E. M., Sippel, S., & Knutti, R. (2021). Increasing probability of record-shattering climate extremes. *Nature Climate Change*, *11*(8), 689–695. <https://doi.org/10.1038/s41558-021-01092-9>
- Friend, A. D., Lucht, W., Rademacher, T. T., Keribin, R., Betts, R., Cadule, P., et al. (2014). Carbon residence time dominates uncertainty in terrestrial vegetation responses to future climate and atmospheric CO₂. *Proceedings of the National Academy of Sciences*, *111*(9), 3280–3285. <https://doi.org/10.1073/pnas.1222477110>
- Fulton, D. J., Clarke, B. J., & Hegerl, G. C. (2023). Bias correcting climate model simulations using unpaired image-to-image translation networks. *Artificial Intelligence for the Earth Systems*, *2*(2), e220031. <https://doi.org/10.1175/AIES-D-22-0031.1>
- Geirinhas, J. L., Russo, A., Libonati, R., Sousa, P. M., Miralles, D. G., & Trigo, R. M. (2021). Recent increasing frequency of compound summer drought and heatwaves in Southeast Brazil. *Environmental Research Letters*, *16*(3), 034036. <https://doi.org/10.1088/1748-9326/abe0eb>
- Golaz, J.-C., Larson, V. E., & Cotton, W. R. (2002). A PDF-based model for boundary layer clouds. Part I: Method and model description. *Journal of the Atmospheric Sciences*, *59*(24), 3540–3551. [https://doi.org/10.1175/1520-0469\(2002\)059<3540:apbmbf>2.0.co;2](https://doi.org/10.1175/1520-0469(2002)059<3540:apbmbf>2.0.co;2)
- Golaz, J.-C., Van Roekel, L. P., Zheng, X., Roberts, A. F., Wolfe, J. D., Lin, W., et al. (2022). The DOE E3SM model version 2: Overview of the physical model and initial model evaluation. *Journal of Advances in Modeling Earth Systems*, *14*(12), e2022MS003156. <https://doi.org/10.1029/2022ms003156>
- Guibas, J., Mardani, M., Li, Z., Tao, A., Anandkumar, A., & Catanzaro, B. (2022). Adaptive Fourier neural operators: Efficient token mixers for transformers. arXiv. <https://doi.org/10.48550/arXiv.2111.13587>
- Harlim, J., Jiang, S. W., Liang, S., & Yang, H. (2021). Machine learning for prediction with missing dynamics. *Journal of Computational Physics*, *428*, 109922. <https://doi.org/10.1016/j.jcp.2020.109922>
- Hausner, M., Orth, R., & Seneviratne, S. I. (2016). Role of soil moisture versus recent climate change for the 2010 heat wave in western Russia. *Geophysical Research Letters*, *43*(6), 2819–2826. <https://doi.org/10.1002/2016GL068036>
- Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., et al. (2020). The ERA5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, *146*(730), 1999–2049. <https://doi.org/10.1002/qj.3803>
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, *9*(8), 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- Holloway, C. E., & Neelin, J. D. (2009). Moisture vertical structure, column water vapor, and tropical deep convection. *Journal of the Atmospheric Sciences*, *66*(6), 1665–1683. <https://doi.org/10.1175/2008JAS2806.1>
- Houser, T., Hsiang, S., Kopp, R., Larsen, K., Delgado, M., Jina, A., et al. (2015). *Economic risks of climate change: An American prospectus*. Columbia University Press. <https://doi.org/10.7312/hous17456>
- Huang, Z., Zhong, L., Ma, Y., & Fu, Y. (2021). Development and evaluation of spectral nudging strategy for the simulation of summer precipitation over the Tibetan Plateau using WRF (v4.0). *Geoscientific Model Development*, *14*(5), 2827–2841. <https://doi.org/10.5194/gmd-14-2827-2021>
- Li, Z., Kovachki, N., Azizzadenesheli, K., Liu, B., Bhattacharya, K., Stuart, A., & Anandkumar, A. (2021). Fourier neural operator for parametric partial differential equations. arXiv. <https://doi.org/10.48550/arXiv.2010.08895>
- Li, Z., Peng, W., Yuan, Z., & Wang, J. (2022). Fourier neural operator approach to large eddy simulation of three-dimensional turbulence. *Theoretical and Applied Mechanics Letters*, *12*(6), 100389. <https://doi.org/10.1016/j.taml.2022.100389>
- Liu, X., Ma, P.-L., Wang, H., Tilmes, S., Singh, B., Easter, R., et al. (2016). Description and evaluation of a new four-mode version of the Modal Aerosol Module (MAM4) within version 5.3 of the Community Atmosphere Model. *Geoscientific Model Development*, *9*(2), 505–522. <https://doi.org/10.5194/gmd-9-505-2016>

- Lucarini, V., Faranda, D., Freitas, A., Freitas, J., Holland, M., Kuna, T., et al. (2016). Extremes and recurrence in dynamical systems (p. 295). <https://doi.org/10.1002/9781118632321>
- Manabe, S., Smagorinsky, J., & Strickler, R. F. (1965). Simulated climatology of a general circulation model with a hydrologic cycle. *Monthly Weather Review*, 93(12), 769–798. [https://doi.org/10.1175/1520-0493\(1965\)093<0769:SCOAGC>2.3.CO;2](https://doi.org/10.1175/1520-0493(1965)093<0769:SCOAGC>2.3.CO;2)
- McGibbon, J., Brenowitz, N. D., Cheeseman, M., Clark, S. K., Dahm, J. P. S., Davis, E. C., et al. (2021). fv3gfs-wrapper: A Python wrapper of the FV3GFS atmospheric model. *Geoscientific Model Development*, 14(7), 4401–4409. <https://doi.org/10.5194/gmd-14-4401-2021>
- McGibbon, J. J., Clark, S. K., Henn, B., Kwa, A., Watt-Meyer, O., Perkins, W. A., & Bretherton, C. S. (2023). Global precipitation correction across a range of climates using CycleGAN (preprint). *Preprints*. <https://doi.org/10.22541/essoar.168881853.36817507/v1>
- Miguez-Macho, G., Stenchikov, G. L., & Robock, A. (2005). Regional climate simulations over North America: Interaction of local processes with improved large-scale flow. *Journal of Climate*, 18(8), 1227–1246. <https://doi.org/10.1175/JCLI3369.1>
- Mintz, Y. (1968). Very long-term global integration of the primitive equations of atmospheric motion: An experiment in climate simulation. In D. E. Billings, J. M. Mitchell, & American Meteorological Society (Eds.), *Causes of Climatic Change: A collection of papers derived from the INQUA—NCAR Symposium on Causes of Climatic Change, August 30–31, 1965, Boulder, Colorado* (pp. 20–36). American Meteorological Society. https://doi.org/10.1007/978-1-935704-38-6_3
- Mlawer, E. J., Taubman, S. J., Brown, P. D., Iacono, M. J., & Clough, S. A. (1997). Radiative transfer for inhomogeneous atmospheres: RRTM, a validated correlated-k model for the longwave. *Journal of Geophysical Research*, 102(D14), 16663–16682. <https://doi.org/10.1029/97jd00237>
- Mons, V., Chassaing, J. C., Gomez, T., & Sagaut, P. (2016). Reconstruction of unsteady viscous flows using data assimilation schemes. *Journal of Computational Physics*, 316, 255–280. <https://doi.org/10.1016/j.jcp.2016.04.022>
- Morrison, H., & Gettelman, A. (2008). A new two-moment bulk stratiform cloud microphysics scheme in the Community Atmosphere Model, version 3 (CAM3). Part I: Description and numerical tests. *Journal of Climate*, 21(15), 3642–3659. <https://doi.org/10.1175/2008jcli2105.1>
- Oleson, K., Lawrence, D., Bonan, G., Drewniack, B., Huang, M., Koven, C., et al. (2013). *Technical description of version 4.5 of the Community Land Model (CLM)* (Technical Note No. NCAR/TN-503+ STR). National Center for Atmospheric Research Earth System Laboratory.
- Pathak, J., Subramanian, S., Harrington, P., Raja, S., Chattopadhyay, A., Mardani, M., et al. (2022). FourCastNet: A Global Data-driven High-resolution Weather Model using Adaptive Fourier Neural Operators. arXiv. <https://doi.org/10.48550/arXiv.2202.11214>
- Payne, A. E., Demory, M.-E., Leung, L. R., Ramos, A. M., Shields, C. A., Rutz, J. J., et al. (2020). Responses and impacts of atmospheric rivers to climate change. *Nature Reviews Earth & Environment*, 1(3), 143–157. <https://doi.org/10.1038/s43017-020-0030-5>
- Qi, D., & Majda, A. J. (2018). Predicting extreme events for passive scalar turbulence in two-layer baroclinic flows through reduced-order stochastic models. *Communications in Mathematical Sciences*, 16(1), 17–51. <https://doi.org/10.4310/CMS.2018.v16.n1.a2>
- Rasp, S., Pritchard, M. S., & Gentine, P. (2018). Deep learning to represent subgrid processes in climate models. *Proceedings of the National Academy of Sciences*, 115(39), 9684–9689. <https://doi.org/10.1073/pnas.1810286115>
- Raymond, C., Horton, R. M., Zscheischler, J., Martius, O., AghaKouchak, A., Balch, J., et al. (2020). Understanding and managing connected extreme events. *Nature Climate Change*, 10(7), 611–621. <https://doi.org/10.1038/s41558-020-0790-4>
- Reynolds, R. W., Rayner, N. A., Smith, T. M., Stokes, D. C., & Wang, W. (2002). An improved in situ and satellite SST analysis for climate. *Journal of Climate*, 15(13), 1609–1625. [https://doi.org/10.1175/1520-0442\(2002\)015<1609:AHSAS>2.0.CO;2](https://doi.org/10.1175/1520-0442(2002)015<1609:AHSAS>2.0.CO;2)
- Robinson, A., Lehmann, J., Barriopedro, D., Rahmstorf, S., & Coumou, D. (2021). Increasing heat and rainfall extremes now far outside the historical climate. *npj Climate and Atmospheric Science*, 4(1), 1–4. <https://doi.org/10.1038/s41612-021-00202-w>
- Sapsis, T. P. (2021). Statistics of extreme events in fluid flows and waves. *Annual Review of Fluid Mechanics*, 53(1), 85–111. <https://doi.org/10.1146/annurev-fluid-030420-032810>
- Schneider, T., Behera, S., Boccaletti, G., Deser, C., Emanuel, K., Ferrari, R., et al. (2023). Harnessing AI and computing to advance climate modelling and prediction. *Nature Climate Change*, 13(9), 887–889. <https://doi.org/10.1038/s41558-023-01769-3>
- Schneider, T., Lan, S., Stuart, A., & Teixeira, J. (2017). Earth system modeling 2.0: A blueprint for models that learn from observations and targeted high-resolution simulations. *Geophysical Research Letters*, 44(24), 12396–12417. <https://doi.org/10.1002/2017GL076101>
- Schneider, T., Teixeira, J., Bretherton, C. S., Brient, F., Pressel, K. G., Schär, C., & Siebesma, A. P. (2017). Climate goals and computing the future of clouds. *Nature Climate Change*, 7(1), 3–5. <https://doi.org/10.1038/nclimate3190>
- Slingo, J., Bates, P., Bauer, P., Belcher, S., Palmer, T., Stephens, G., et al. (2022). Ambitious partnership needed for reliable climate prediction. *Nature Climate Change*, 12(6), 499–503. <https://doi.org/10.1038/s41558-022-01384-8>
- Smagorinsky, J. (1963). General circulation experiments with the primitive equations: I. The basic experiment. *Monthly Weather Review*, 91(3), 99–164. [https://doi.org/10.1175/1520-0493\(1963\)091<0099:GCEWTP>2.3.CO;2](https://doi.org/10.1175/1520-0493(1963)091<0099:GCEWTP>2.3.CO;2)
- Smagorinsky, J., Manabe, S., & Holloway, J. L. (1965). Numerical results from a nine-level general circulation model of the atmosphere. *Monthly Weather Review*, 93(12), 727–768. [https://doi.org/10.1175/1520-0493\(1965\)093<0727:NRFANL>2.3.CO;2](https://doi.org/10.1175/1520-0493(1965)093<0727:NRFANL>2.3.CO;2)
- Stensrud, D. J. (2007). *Parameterization schemes: Keys to understanding numerical weather prediction models*. Cambridge University Press. <https://doi.org/10.1017/CBO9780511812590>
- Stevens, B., Satoh, M., Auger, L., Biercamp, J., Bretherton, C. S., Chen, X., et al. (2019). DYAMOND: The Dynamics of the Atmospheric general circulation Modeled On Non-hydrostatic Domains. *Progress in Earth and Planetary Science*, 6(1), 61. <https://doi.org/10.1186/s40645-019-0304-z>
- Storch, H. V., Langenberg, H., & Feser, F. (2000). A spectral nudging technique for dynamical downscaling purposes. *Monthly Weather Review*, 128(10), 3664–3673. [https://doi.org/10.1175/1520-0493\(2000\)128<3664:ASNTFD>2.0.CO;2](https://doi.org/10.1175/1520-0493(2000)128<3664:ASNTFD>2.0.CO;2)
- Sun, J., Zhang, K., Wan, H., Ma, P.-L., Tang, Q., & Zhang, S. (2019). Impact of nudging strategy on the climate representativeness and hindcast skill of constrained EAMv1 simulations. *Journal of Advances in Modeling Earth Systems*, 11(12), 3911–3933. <https://doi.org/10.1029/2019MS001831>
- Takens, F. (1981). Detecting strange attractors in turbulence. In *Dynamical Systems and Turbulence, Warwick 1980: Proceedings of a symposium held at the University of Warwick 1979/80* (pp. 366–381). Springer.
- Taylor, M. A., Cyr, A. S., & Fournier, A. (2009). A non-oscillatory advection operator for the compatible spectral element method. In *International Conference on Computational Science* (pp. 273–282).
- Taylor, M. A., Guba, O., Steyer, A., Ullrich, P. A., Hall, D. M., & Eldred, C. (2020). An energy consistent discretization of the nonhydrostatic equations in primitive variables. *Journal of Advances in Modeling Earth Systems*, 12(1), e2019MS001783. <https://doi.org/10.1029/2019ms001783>
- Tomita, H., Miura, H., Iga, S., Nasuno, T., & Satoh, M. (2005). A global cloud-resolving simulation: Preliminary results from an aqua planet experiment. *Geophysical Research Letters*, 32(8), L08805. <https://doi.org/10.1029/2005GL022459>
- Vlachas, P. R., Byeon, W., Wan, Z. Y., Sapsis, T. P., & Koumoutsakos, P. (2018). Data-driven forecasting of high-dimensional chaotic systems with long short-term memory networks. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 474(2213), 20170844. <https://doi.org/10.1098/rspa.2017.0844>

- Wan, Z. Y., Vlachas, P., Koumoutsakos, P., & Sapsis, T. (2018). Data-assisted reduced-order modeling of extreme events in complex dynamical systems. *PLoS One*, *13*(5), e0197704. <https://doi.org/10.1371/journal.pone.0197704>
- Watt-Meyer, O., Brenowitz, N. D., Clark, S. K., Henn, B., Kwa, A., McGibbon, J., et al. (2021). Correcting weather and climate models by machine learning nudged historical simulations. *Geophysical Research Letters*, *48*(15), e2021GL092555. <https://doi.org/10.1029/2021GL092555>
- Wedi, N. P., Polichtchouk, I., Dueben, P., Anantharaj, V. G., Bauer, P., Boussetta, S., et al. (2020). A baseline for global weather and climate simulations at 1 km resolution. *Journal of Advances in Modeling Earth Systems*, *12*(11), e2020MS002192. <https://doi.org/10.1029/2020MS002192>
- Wikner, A., Harvey, J., Girvan, M., Hunt, B. R., Pomerance, A., Antonsen, T., & Ott, E. (2022). Stabilizing machine learning prediction of dynamics: Noise and noise-inspired regularization. arXiv. <https://doi.org/10.48550/arXiv.2211.05262>
- Witte, J. C., Douglass, A. R., da Silva, A., Torres, O., Levy, R., & Duncan, B. N. (2011). NASA A-Train and Terra observations of the 2010 Russian wildfires. *Atmospheric Chemistry and Physics*, *11*(17), 9287–9301. <https://doi.org/10.5194/acp-11-9287-2011>
- Wood, R. (2012). Stratocumulus clouds. *Monthly Weather Review*, *140*(8), 2373–2423. <https://doi.org/10.1175/MWR-D-11-00121.1>
- Yuval, J., & O’Gorman, P. A. (2020). Stable machine-learning parameterization of subgrid processes for climate modeling at a range of resolutions. *Nature Communications*, *11*(1), 3295. <https://doi.org/10.1038/s41467-020-17142-3>
- Yuval, J., O’Gorman, P. A., & Hill, C. N. (2021). Use of neural networks for stable, accurate and physically consistent parameterization of subgrid atmospheric processes with good performance at reduced precision. *Geophysical Research Letters*, *48*(6), e2020GL091363. <https://doi.org/10.1029/2020GL091363>
- Zhang, G. J., & McFarlane, N. A. (1995). Sensitivity of climate simulations to the parameterization of cumulus convection in the Canadian climate centre general circulation model. *Atmosphere-Ocean*, *33*(3), 407–446. <https://doi.org/10.1080/07055900.1995.9649539>
- Zhang, H., Harlim, J., & Li, X. (2021). Error bounds of the invariant statistics in machine learning of ergodic Itô diffusions. *Physica D: Nonlinear Phenomena*, *427*, 133022. <https://doi.org/10.1016/j.physd.2021.133022>
- Zhang, S., Harrop, B., Leung, L., Charalampopoulos, A.-T., Barthel, B., Xu, W., & Sapsis, T. (2023). A machine learning bias correction of large-scale environment of extreme weather events in E3SM atmosphere model. *Authorea Preprints*. <https://doi.org/10.22541/essoar.170067232.22392274/v1>
- Zhang, S., Zhang, K., Wan, H., & Sun, J. (2022). Further improvement and evaluation of nudging in the E3SM Atmosphere Model version 1 (EAMv1): Simulations of the mean climate, weather events, and anthropogenic aerosol effects. *Geoscientific Model Development*, *15*(17), 6787–6816. <https://doi.org/10.5194/gmd-15-6787-2022>
- Zscheischler, J., Westra, S., van den Hurk, B. J. J. M., Seneviratne, S. I., Ward, P. J., Pitman, A., et al. (2018). Future climate risk from compound events. *Nature Climate Change*, *8*(6), 469–477. <https://doi.org/10.1038/s41558-018-0156-3>