



## RESEARCH ARTICLE

10.1029/2018MS001472

## Applications of Deep Learning to Ocean Data Inference and Subgrid Parameterization

## Key Points:

- We successfully use convolutional neural networks to predict unresolved turbulent processes and subsurface velocities
- The neural networks generalize to different regions, dynamical regimes, and forcing
- Global momentum conservation for eddy parameterization can be respected without sacrificing accuracy

## Correspondence to:

T. Bolton,  
tom.bolton@physics.ox.ac.uk

## Citation:

Bolton, T., & Zanna, L. (2019). Applications of deep learning to ocean data inference and subgrid parameterization. *Journal of Advances in Modeling Earth Systems*, 11, 376–399. <https://doi.org/10.1029/2018MS001472>

Received 13 AUG 2018

Accepted 27 DEC 2018

Accepted article online 4 JAN 2019

Published online 30 JAN 2019

Corrected 22 FEB 2019

This article was corrected on 22 FEB 2019. See the end of the full text for details.

Thomas Bolton<sup>1</sup> and Laure Zanna<sup>1</sup> <sup>1</sup>Department of Physics, University of Oxford, Oxford, UK

**Abstract** Oceanographic observations are limited by sampling rates, while ocean models are limited by finite resolution and high viscosity and diffusion coefficients. Therefore, both data from observations and ocean models lack information at small and fast scales. Methods are needed to either extract information, extrapolate, or upscale existing oceanographic data sets, to account for or represent unresolved physical processes. Here we use machine learning to leverage observations and model data by predicting unresolved turbulent processes and subsurface flow fields. As a proof of concept, we train convolutional neural networks on degraded data from a high-resolution quasi-geostrophic ocean model. We demonstrate that convolutional neural networks successfully replicate the spatiotemporal variability of the subgrid eddy momentum forcing, are capable of generalizing to a range of dynamical behaviors, and can be forced to respect global momentum conservation. The training data of our convolutional neural networks can be subsampled to 10–20% of the original size without a significant decrease in accuracy. We also show that the subsurface flow field can be predicted using only information at the surface (e.g., using only satellite altimetry data). Our results indicate that data-driven approaches can be exploited to predict both subgrid and large-scale processes, while respecting physical principles, even when data are limited to a particular region or external forcing. Our in-depth study presents evidence for the successful design of ocean eddy parameterizations for implementation in coarse-resolution climate models.

**Plain Language Summary** Models of the ocean and ocean observations are imperfect. Due to this imperfection, simulations of the ocean and our observations are not quite the same as the true ocean currents. We, therefore, need ways to make our ocean data more realistic and complete and to make it more similar to the actual ocean. Scientists have traditionally approached this problem in a pen-and-paper style, considering physical theories and mechanisms. This study instead uses machine learning, which focuses on data as opposed to equations on a black board. We successfully use a particular type of machine learning algorithm, called a convolutional neural network, to make the most of current oceanographic data. This type of neural network works well even if ocean data are limited to a particular area. Future work will involve combining machine learning with physical theories of the ocean.

## 1. Introduction

Satellite observations have produced a wealth of information on the ocean circulation (Abernathey & Marshall, 2013; Chelton et al., 2007; Greatbatch et al., 2010; Le Traon & Morrow, 2001; Morrow et al., 1994; Scott & Wang, 2005). However, raw satellite altimetry data subsamples the ocean and does not measure subsurface quantities. Temporally, measurements at the same location are made twice every orbital cycle, while the spatial sampling depends upon the distance between ground tracks. To improve the subsampling rates, measurements from multiple satellites are combined (Le Traon et al., 1998) to produce an optimal estimate.

The process of combining measurements from multiple satellites includes spatiotemporal filtering, which leads to a more *smoothed* view of the dynamical processes at the oceans surface, removing variability due to mesoscale and submesoscale eddies. The filtering can also lead to spurious physical signals, as studied by Arbic et al. (2013), which showed that filtering data can lead to exaggerated forward cascades of energy. The new Surface Water and Ocean Topography mission will have a large swath of 120 km, providing unprecedented detail on the oceans surface. Despite the high spatial sampling rate, measurements may still be limited by the temporal sampling rate of 11 days (Durand et al., 2010).

©2019. The Authors.

This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

Similar to satellite observations, Ocean General Circulation Models (OGCMs) are useful for studying ocean dynamics. However, high-resolution models are computationally expensive, and the current resolution of models is not high enough to fully resolve the first baroclinic deformation radius at midlatitudes (Hallberg, 2013). Also, due to their finite resolution, they require large viscosity and diffusion coefficients in order to remain numerically stable (Jochum et al., 2008). The combination of finite resolution and artificially high viscosity, diffuses momentum and smooths out features such as jets and mesoscale eddies (Hewitt et al., 2016; Kjellsson & Zanna, 2017).

Therefore, both observations and models are missing the interactions of oceanic turbulence at small scales, which play an important role in maintaining the large-scale circulation (Greatbatch, Zhai, Claus, et al., 2010; Greatbatch, Zhai, Kohlmann, & Czeschel, 2010; Kang & Curchitser, 2015; Waterman & Jayne, 2010; Waterman et al., 2011); with satellite observations only providing surface information. We thus consider the general problem: given some smoothed view of the oceans surface, what information can be generated on small-scale turbulent interactions and subsurface quantities? Illuminating unresolved quantities using *seen* quantities would extend the reach of existing data sets and could potentially improve the representations of unresolved eddies in OGCMs.

We tackle this problem with machine learning. Machine learning has grown in popularity in recent years and has been applied to weather prediction (Esteves et al., 2018; McGovern et al., 2017), climate model parameter sensitivity studies (Anderson & Lucas, 2018), chaotic dynamical systems forecasting (Pathak, Hunt, et al., 2018; Pathak, Wikner, et al., 2018; Vlachas et al., 2018), and parameterizing unresolved atmospheric processes (Brenowitz & Bretherton, 2018; Gentine et al., 2018; Jiang et al., 2018; O'Gorman & Dwyer, 2018). The foundational principle of machine learning is extracting information from data. When used to improve our understanding of the Earth system, these data-driven methods are an empirical bottom-up approach, whereas the rationalist top-down approach considers physical principles and mechanisms. Here we take the empirical route by exploiting recent developments in machine learning.

Using empirical methods to leverage ocean observations is not new. For example, using satellite altimetry data, Keating et al. (2012) constructed a stochastic model to *super-resolve* the velocity field and predict the velocity at depth. Similarly, Keating and Smith (2015) used a stochastic model to produce a super-resolved sea surface temperature (SST) field, given a low-resolution observation of SST. With regard to machine learning, Chapman and Charantonis (2017) constructed a form of neural network known as a self-organizing map to reconstruct subsurface velocities in the Southern Ocean using satellite altimetry data and Argo floats. Other studies have used random forests to predict subsurface temperature anomalies (Su et al., 2018) and Southern Ocean oxygen content (Giglio et al., 2018).

In the previous studies that leverage oceanic observations, there is an abundance of coarse-resolution data (satellite altimetry) but limited data on the desired quantities (e.g., high-resolution SST or Argo subsurface velocities); as is the case with OGCMs, where high-resolution data are less readily available due to the computational cost. A similar challenge is when data are only available for particular regions, such as mooring data (Hogg, 1992) or gliders (Davis et al., 2008; Rudnick et al., 2004). A machine learning algorithm trained on region-limited data would have to adapt to new regions with different physics; this task is well suited to deep neural networks, which are known for a strong ability to generalize (Goodfellow et al., 2016; Krizhevsky et al., 2012; LeCun et al., 2015).

However, deep neural networks are typically considered a *black box*; that is, they lack simple interpretations. It is therefore difficult to assess whether such data-driven methods respect physical principles (e.g., conservation of energy or momentum). For example, neural networks have been used to develop Reynolds-averaged turbulence models (Kutz, 2017; Tracey et al., 2015), where the studies of Ling, Jones, and Templeton (2016) and Ling, Kurzawski, and Templeton (2016), in particular, show that a neural network can respect Galilean invariance by utilizing the invariant tensors of Pope (1975). The studies of Ling, Jones, and Templeton (2016) and Ling, Kurzawski, and Templeton (2016) are important in moving toward data-driven approaches that respect the physical properties of the system.

In this paper we focus on a particular machine learning algorithm, namely, convolutional neural networks, in order to leverage observations and coarse-resolution model data. Our aim is to test whether they can be used to reveal information on unresolved turbulent processes and subsurface flow fields, and to determine

if they are suited to situations where data are limited to a particular region. To move toward these aims, as a proof of concept, we will address the following questions:

1. Can convolutional neural networks represent the spatiotemporal variability of the subgrid eddy momentum forcing?
2. How sensitive are the neural networks to the physical processes occurring within each region, and how well do they generalize to ocean models in different configurations?
3. Is it possible to physically constrain neural networks to respect global momentum conservation?
4. Using only information at the surface, can neural networks predict the subsurface flow fields?

By using data from an idealized high-resolution ocean model, we show that convolutional neural networks can represent both the spatial and temporal variability of the eddy momentum forcing. The region the neural network is trained on and therefore the dynamical processes occurring within that region significantly impact the performance of the neural network. In particular, training on the most turbulent region produces the best overall performing neural network. The neural networks successfully generalize to models with different viscosity coefficients and external wind forcings. Initially, momentum is not conserved globally, but the neural networks can be constrained to respect momentum conservation without a significant reduction in accuracy. A neural network can accurately predict the subsurface flow field when there is a strong barotropic component to the flow.

The paper is organized as follows. The quasi-geostrophic ocean model, the degrading of model data, and convolutional neural network are introduced in section 2. Performance diagnostics of the neural networks, in terms of nonlocal predictions and generalizing to different model configurations, are presented in section 3. We explore methods of physically constraining the neural networks in section 4. Section 5 presents a neural network trained to predict subsurface flow fields using only information at the surface. We summarize and discuss our results in section 6.

## 2. Data and Methods

### 2.1. Quasi-Geostrophic Ocean Model

We use the PEQUOD model which solves the three-dimensional baroclinic quasi-geostrophic (QG) potential vorticity equation, with constant wind forcing on a beta plane (e.g., Berloff, 2005). The model has a bounded-square domain with a flat bottom.

The configuration of this model leads to two large-scale circulation gyres separated latitudinally by a strong meandering zonal jet. The model is configured to represent an idealized version of current systems such as the Gulf Stream in the North Atlantic or the Kuroshio Extension in the North Pacific; both these current systems exhibit vigorous eddies interacting with a strong mean flow. The time mean stream function, which illustrates the double-gyre flow structure, can be seen in Figure 1a of Mana and Zanna (2014).

The potential vorticity  $q$  is given by

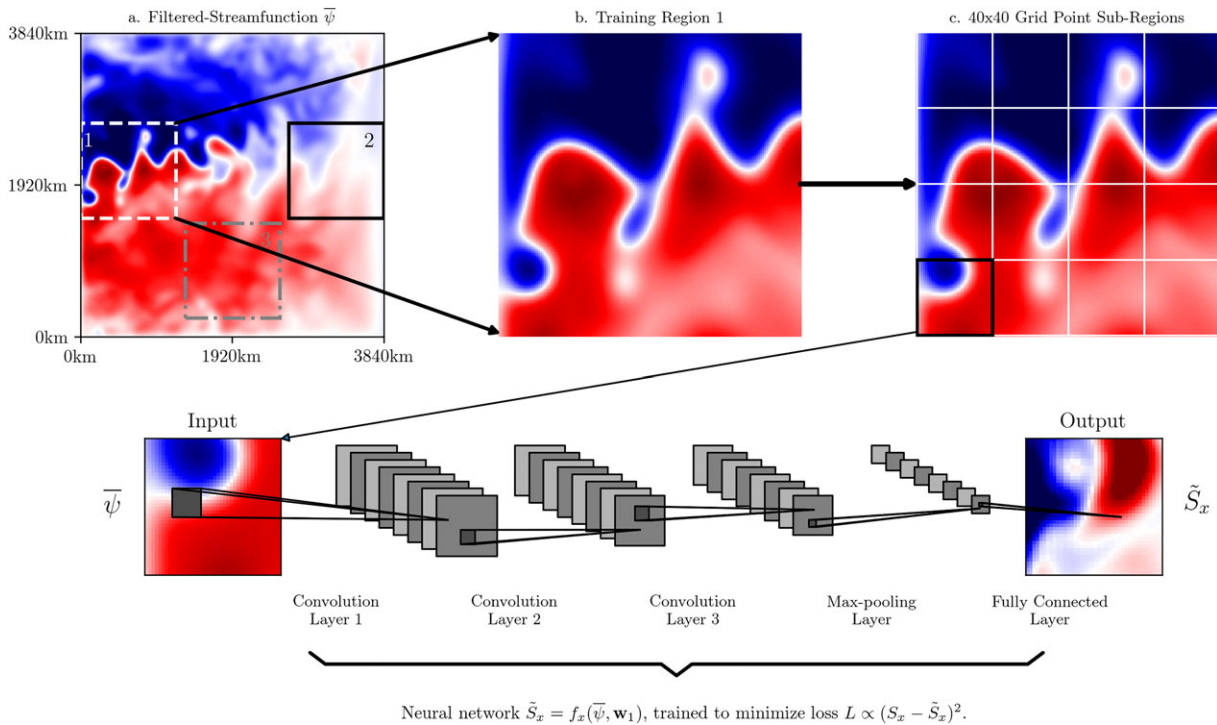
$$q = \nabla^2 q + \beta y + \frac{\partial}{\partial z} \left( \frac{f_0^2}{N^2} \frac{\partial \psi}{\partial z} \right), \quad (1)$$

where  $f = f_0 + \beta y$  is the planetary vorticity,  $f_0$  is the Coriolis parameter,  $\beta = df/dy$  is the Rossby parameter,  $\nabla = (\partial/\partial x, \partial/\partial y)$  is the horizontal gradient operator,  $N = \left(-\frac{g}{\rho} \frac{d\rho}{dz}\right)^{1/2}$  is the Brunt-Väisälä frequency,  $g$  is gravity,  $\rho$  is density, and  $\psi$  is the stream function for the nondivergent horizontal velocity  $\mathbf{u} = (-\partial\psi/\partial y, \partial\psi/\partial x)$ .

The model has three layers ( $m = 1$  upper,  $m = 2$  middle, and  $m = 3$  upper), with thicknesses  $H_m$  of 250, 750, and 3,000 m, respectively. For each layer, the following prognostic equation is solved

$$\frac{\partial q}{\partial t} + (\mathbf{u} \cdot \nabla)q = D + \mathcal{F}, \quad (2)$$

where  $D = \nu \nabla^4 \psi - r \nabla^2 \psi \delta_{m,3}$  is the dissipation, and  $\mathcal{F} = (\nabla \times \tau)_z \delta_{m,1} / \rho_0 H_1$  is the applied wind stress curl forcing, where  $\delta_{i,j}$  is the Kronecker delta function. The horizontal resolution of the model is 7.5 km, such that the model is eddy resolving. The first term in the dissipation is a fourth-order term equivalent to Laplacian viscosity, with viscosity coefficient  $\nu$ . The second dissipation term parameterizes the presence of an Ekman



**Figure 1.** Panel (a) illustrates the upper-layer filtered-stream function  $\bar{\psi}$  of the quasi-geostrophic model, including the three regions in which we train the neural networks: region 1 (white dashed) is on the western boundary, region 2 (black solid) is on the eastern boundary, and region 3 (gray dash dotted) is centered on the southern gyre. Panel (b) shows a close-up of the filtered-stream function  $\bar{\psi}$  within training region 1 while panel (c) illustrates how training region 1 is split into sixteen  $40 \times 40$  grid point subregions—the size of the input and output arrays of the neural network is  $40 \times 40$  grid points. The input variable of each neural network is the filtered-stream function  $\bar{\psi}$ , and the output variable is either the zonal component  $\tilde{S}_x$  or meridional component  $\tilde{S}_y$  of the subfilter eddy momentum forcing. The architecture of the convolutional neural network, with an example input  $\bar{\psi}$  and output  $\tilde{S}_x$ , is illustrated underneath panels (a)–(c).

layer with bottom drag coefficient  $r$  (and therefore only acts on the bottom  $m = 3$  layer). The wind stress forcing applied to the upper  $m = 1$  layer is given explicitly by

$$F(x, y) = \begin{cases} -\tau_0 \frac{0.92\pi}{L\rho_0 H_1} \sin\left(\frac{\pi y}{g(x)}\right) & y \leq g(x), \\ \tau_0 \frac{2\pi}{0.9L\rho_0 H_1} \sin\left(\frac{\pi[2y-g(x)]}{L-g(x)}\right) & y > g(x), \end{cases} \quad (3)$$

where  $g(x) = L/2 + 0.2(x - L/2)$ ,  $L = 3,840$  km is the domain length, and  $\rho_0$  is the reference density. After the model has been integrated from rest to a statistically steady state, we save 10 years of model output at daily resolution of the turbulent double-gyre circulation. For further details on the QG model, see Mana and Zanna (2014) and Zanna et al. (2017), and for a list of the model parameters, see Table 1. We use the data generated by the ocean model to train various neural networks, but only after degrading the data, to make it similar to observations or low-resolution model.

## 2.2. Degrading High-Resolution Data

We degrade the fields from the high-resolution QG model using a spatial 2-D low-pass filter, in order to produce data that are similar to satellite altimetry or a model with a large numerical dissipation. From the filtering of the model data, we can then calculate the forcing from unresolved small-scale turbulent processes.

At every time slice in the data, we take a high-resolution variable  $a$  at a particular layer and apply a two-dimensional spatial Gaussian filter. We denote filtered variables as  $\bar{a}$ , and subfilter variables as the deviation from the filtered variable  $a' = a - \bar{a}$ . The value of a function  $a(x, y)$ , after the Gaussian low-pass filtering

**Table 1**

*Details on the Following: The Quasi-Geostrophic Ocean Model Parameters, the Data Sets Used to Train the Neural Networks, the Architecture Parameters, and the Optimization Parameters*

Quasi-Geostrophic Model Parameters	
Domain size (grid points)	512 × 512
Domain length ( $L$ )	3,840 km
Resolution ( $\Delta x$ )	7.5 km
Viscosity ( $\nu$ )	75 m <sup>2</sup> /s
Rossby deformation radii ( $L_{Ro}$ )	40, 23 km
Velocity scale ( $\sqrt{EKE}$ )	0.21 m/s
Planetary vorticity ( $f_0$ )	10 <sup>-4</sup> s <sup>-1</sup>
Rossby parameter ( $\beta$ )	2 × 10 <sup>-11</sup> m <sup>-1</sup> /s
Gravity ( $g$ )	9.8 m/s <sup>2</sup>
Reduced gravity ( $g'$ )	0.034, 0.018 m/s <sup>2</sup>
Bottom drag coefficient ( $r$ )	4 × 10 <sup>-8</sup> s <sup>-1</sup>
Wind stress amplitude ( $\tau_0$ )	0.8 N/m <sup>2</sup>
Reference density ( $\rho_0$ )	10 <sup>3</sup> kg/m <sup>3</sup>
Neural Network Data Details	
Data source	Quasi-geostrophic ocean model
Input variable (feature)	Filtered-stream function $\bar{\psi}$
Output variables (targets)	Subfilter momentum forcing $S_x, S_y$
Training region 1	Western boundary
Training region 2	Eastern boundary
Training region 3	Southern gyre
Number of training samples	5,800 (years 1–9)
Number of validation samples	5,600 (year 10)
Standardization method	Zero mean, unit variance
Neural Network Architecture	
Input size	40 × 40
Number of convolution layers	3
Number of filters for each convolution layer	16, 16*8, 8*8
Size of filter for each convolution layer	8 × 8, 4 × 4, 4 × 4
Filter stride for each convolution layer	2, 1, 1
Activation function for each convolution layer	SELU, SELU, SELU
Max pooling kernel size	2
Output layer activation function	None/Linear
Output size	40 × 40
Neural Network Training Parameters	
Loss function	Mean-square error
Optimizer	Adam
Learning rate	0.001
Momentum	0.9
Batch size	16
Training epochs	200

operation  $G \star a$  at a point  $(x_0, y_0)$ , is given by

$$\begin{aligned} \bar{a}(x_0, y_0) &= G \star a = \iint a(x, y)G(x_0, y_0, x, y)dx dy \\ &= \frac{1}{2\pi\sigma^2} \iint a(x, y)e^{-((x-x_0)^2+(y-y_0)^2)/2\sigma^2} dx dy, \end{aligned} \quad (4)$$

where  $\sigma = 30$  km is the standard deviation of the Gaussian filter, which determines the length scale at which information (below that length scale) is removed. Therefore, the filter acts to remove information on dynamical processes at spatial scales smaller than 30 km.

Using the low-pass filter defined in equation (4), we can now express the effects of the unresolved (subfilter) variables onto the resolved (filtered) variables. Ignoring vertical effects and planetary vorticity, the horizontal momentum equation is given by

$$\frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} = \mathbf{F} + \mathbf{D}, \quad (5)$$

where  $\mathbf{F}$  and  $\mathbf{D}$  are the momentum forcing and dissipation, respectively. Applying a low-pass filter to equation (5), and then adding  $(\bar{\mathbf{u}} \cdot \nabla) \bar{\mathbf{u}}$  to both sides of the equation, leads to

$$\frac{\partial \bar{\mathbf{u}}}{\partial t} + (\bar{\mathbf{u}} \cdot \nabla) \bar{\mathbf{u}} = \bar{\mathbf{F}} + \bar{\mathbf{D}} + \left[ (\bar{\mathbf{u}} \cdot \nabla) \bar{\mathbf{u}} - \overline{(\mathbf{u} \cdot \nabla) \mathbf{u}} \right], \quad (6)$$

$$\frac{\partial \bar{\mathbf{u}}}{\partial t} + (\bar{\mathbf{u}} \cdot \nabla) \bar{\mathbf{u}} = \bar{\mathbf{F}} + \bar{\mathbf{D}} + \mathbf{S}, \quad (7)$$

$$\text{where } \mathbf{S} = \underbrace{(\bar{\mathbf{u}} \cdot \nabla) \bar{\mathbf{u}} - \overline{(\mathbf{u} \cdot \nabla) \mathbf{u}}}_{\text{Subfilter eddy momentum forcing}}. \quad (8)$$

The low-pass filtering operation results in an additional forcing term in equation (7) for the filtered momentum; the additional momentum forcing  $\mathbf{S}$  is given by equation (8), the divergence of a Reynolds stress. The vector  $\mathbf{S} = (S_x, S_y)$  represents the effects of the subfilter momentum field on the filtered momentum field, that is, the interaction between small-scale eddies and the large-scale flow. As the subfilter eddy forcing  $\mathbf{S}$  depends on the subfilter variables, it requires a physical parameterization or closure.

### 2.3. Predictive Algorithm: Convolutional Neural Networks

Convolutional Neural Networks (CNNs) have proven successful in many areas of computer vision (Dong et al., 2016; Krizhevsky et al., 2012; Simonyan & Zisserman, 2014), where the primary objective is to extract information from an image, in order to perform a particular task. CNNs work by applying successive layers of convolutions (a form of spatial filtering) to the input; the complexity of the extracted information increases with the number of convolution layers. The powerful property of CNNs is that the filters of each convolution are learned as part of the training process—they are not specified a priori. Therefore, CNNs learn to extract the most *useful* information from the input variable, given training on a particular data set.

We chose to use CNNs, as opposed to a deep neural network of multiple fully connected layers, due to their superior performance in computer vision tasks where the inputs have a two-dimensional structure (Krizhevsky et al., 2012). We wanted a machine learning algorithm that could exploit the two dimensional lateral structure of turbulent fluids. Spatial filtering of the equations of motion of turbulent fluids is not new and is used in large eddy simulation (Moeng, 1984; Sagaut, 2006). Therefore, the learned filtering operations of a CNN appeared to be a natural choice of data-driven algorithm to apply to geophysical flows.

The training process involves the minimization of an appropriately defined loss function, which measures the difference between the output of the CNN, and the desired targets. If the optimization procedure was successful, such that the loss function on previously unseen data converges, the CNN will have learned to extract the most important information from the input. The CNN then uses the information to predict continuous values. The CNN constructs the final prediction through a linear regression layer, which regresses the desired output onto the final feature maps (feature maps are the intermediate results of each convolution layer).

Here we use CNNs to represent the subfilter eddy momentum forcing. The input is the filtered-stream function  $\bar{\psi}$  of the upper vertical layer, which represents our resolved variable that the neural networks will extract information from. The output variables are the zonal  $S_x$  and meridional  $S_y$  components of the subfilter momentum forcing  $\mathbf{S}$ , defined by equation (8). An example input and output is shown in Figure 1. Separate CNNs are trained for each component of the subfilter momentum forcing  $S_x$  and  $S_y$ . We only consider data from the upper layer of the model; this is because the flow is surface-intensified, and we are assuming that our filtered quantities are similar to satellite altimetry data, which only provide information at the surface.

In addition to testing whether it is possible to train a neural network to predict  $S_x$  and  $S_y$ , from  $\bar{\psi}$ , we explore how a neural network trained on one region performs on another previously unseen region, that is, how important local versus nonlocal information is for different regions. We therefore construct three different data sets from the QG model data, one for each region being studied. We choose regions which differ most in their dynamical behavior, and are shown in Figure 1a: Region 1 is near the jet separation point of the western boundary, where there is a strong, inertial zonal jet. Region 2 is near the eastern boundary downstream of the jet extension, where the dynamics are more wave-like in nature. Region 3 is in the center of the southern gyre, which is energetically less active than regions 1 and 2.

Data from the three regions are split temporally into training and validation data sets. The 10 years of daily data (3,650 days) are split into the first  $\sim 9$  years (3,300 days) to train the neural networks, and the final year (350 days) is set aside for validation. To reduce the computational cost and the number of parameters of each CNN we split each region spatially from the initial  $160 \times 160$  grid points, to sixteen  $40 \times 40$  grid point subregions, as depicted in Figure 1c. Reducing the input and output size of the neural network from  $160 \times 160$  to  $40 \times 40$  significantly decreases the number of trainable weights, and therefore the computational cost (we attempted to make predictions for the full  $160 \times 160$  of each training region, but this led to a neural network with over 250,000,000 parameters, which was computationally impractical).

Making predictions for a  $40 \times 40$  area instead of a  $160 \times 160$  area also increases the amount of training and validation data by a factor of 16, from 3,300 and 350 samples to 52,800 and 5,600, respectively, where a sample is defined as a single input-output pair of the neural network. We therefore have 52,800 spatial maps (size  $40 \times 40$  grid points) of input-output pairs to train the neural networks, and 5,600 spatial maps of input-output pairs set aside for validation.

We train CNNs to separately predict  $S_x$  and  $S_y$ , using data from three different regions of the model; this gives a total of six neural networks. Each neural network is denoted by  $f_i(\bar{\psi}, \mathbf{w}_R)$ , where  $i = (x, y)$  refers to the component of  $\mathbf{S}$  being predicted,  $\mathbf{w}_R$  are the trained weights of the neural network, and  $R = 1, 2, 3$  refers to the region on which the neural network has been trained. For example, the neural network trained on region 2 to predict the meridional component  $S_y$  is denoted by  $f_y(\bar{\psi}, \mathbf{w}_2)$ .

To distinguish predictions from the true values, we label neural network predictions as  $\tilde{S}_x = f_x(\bar{\psi}, \mathbf{w}_R)$ , and  $\tilde{S}_y = f_y(\bar{\psi}, \mathbf{w}_R)$ , while the true values of the subfilter momentum forcing remain as  $S_x, S_y$ . We use the mean-square error as the loss function,

$$L = \sum (S_x - \tilde{S}_x)^2, \text{ or } \sum (S_y - \tilde{S}_y)^2, \quad (9)$$

which quantifies the difference between the neural network predictions and the truth, and where the summation is over all samples. The neural networks are trained (i.e., optimized) using a form of stochastic gradient descent, namely, the Adam optimization algorithm (Kingma & Ba, 2014), which minimizes the loss function  $L$  defined in equation (9). The training of each neural network  $f_i(\bar{\psi}, \mathbf{w}_R)$ , iteratively adjusts the values of the weights  $\mathbf{w}_R$ , such that the loss function in equation (9) is minimized. Therefore, each neural network has a different set of weights  $\mathbf{w}_R$ ; it is these weights which determine how each neural network extracts information and makes predictions.

The architecture used for each  $f_i(\bar{\psi}, \mathbf{w}_R)$  contains three convolution layers, a max pooling layer and a final fully connected layer (Figure 1). The max pooling layer reduces the dimensionality of the previous layer, by selecting the maximum value within a  $2 \times 2$  grid point area—max pooling is effective when there is significant correlation between points in the feature maps. To give the neural networks the ability to learn nonlinear functions, activation functions are added between layers. Here we use the scaled exponential linear unit (SELU; Klambauer et al., 2017). SELU activation functions scale the data toward zero mean and

unit variance, removing the need for batch normalization—batch normalization enforces zero mean and unit variance at each stage of the network but requires additional training.

The specific architecture was constructed by adjusting all parameters and observing which configuration most effectively minimizes the loss function on the validation data. See Table 1 for more details of the architecture and training procedure. The total number of parameters of each neural network is 325,728.

We train and implement each neural network using Keras (Chollet, 2015), with the Tensorflow backend (Abadi et al., 2016). Before training, all data sets are separately normalized to zero mean and unit variance. Each CNN is trained for 200 epochs (1 epoch = 1 full pass of all the training data through the optimization algorithm), taking approximately 10 CPU hours, after which there is negligible change in the loss function of the validation data.

Once all six neural networks are trained, we make the predictions  $\tilde{S}_x$  and  $\tilde{S}_y$  using the filtered-stream function  $\tilde{\psi}$  from the validation data set, that is, the final year of withheld data. We make predictions for the full domain to determine how each neural network generalizes to unseen, dynamically distinct, regions. As the input and output size of each neural network is  $40 \times 40$  grid points, we tile together predictions for the full domain of size  $512 \times 512$ ; the tiling leads to errors at the boundaries of each tile, where discontinuities can emerge. To reduce the tiling error, we make predictions using overlapping tiles and then average the results at each grid point.

In order to make predictions of the subsurface flow field, using only information at the surface, we train a new neural network. The new neural network has an identical architecture to those discussed previously and is trained to predict the middle-layer stream function using the upper-layer stream function as the input; this neural network is described in more detail in section 5.

### 3. Neural Network Generalization and Sensitivity

#### 3.1. Nonlocal Predictions

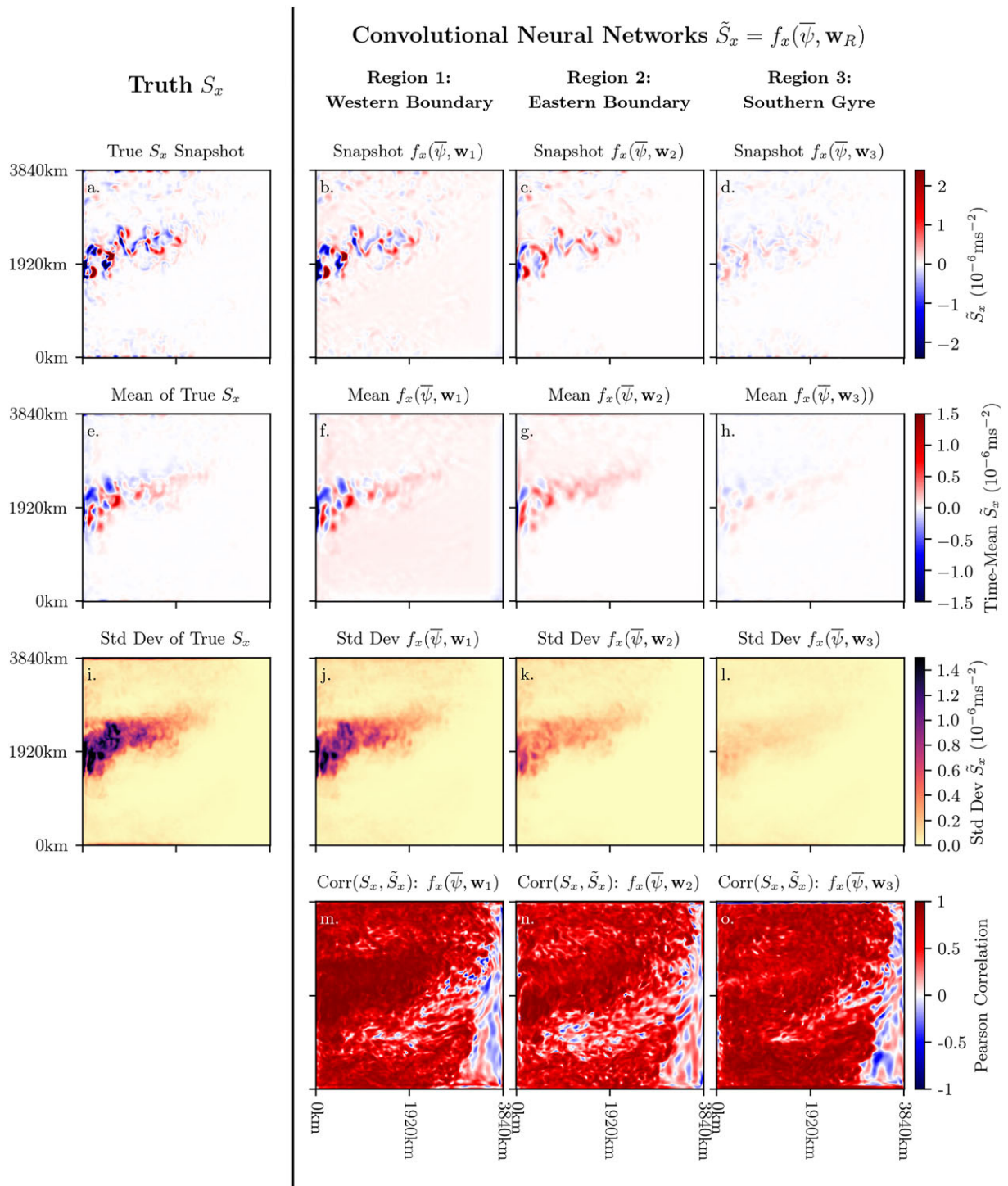
The filtered-stream function represents, for example, observational measurements from satellite altimetry or coarse-resolution model data. The subfilter eddy momentum forcing represents unresolved turbulent processes. Our goal is to replicate the complex spatiotemporal variability of  $S_x$  and  $S_y$  using neural networks  $f_i(\tilde{\psi}, \mathbf{w}_R)$ . However observational data such as moorings (Hogg, 1992) or gliders (Davis et al., 2008; Rudnick et al., 2004), may only be available for a particular region; we therefore only train the neural networks using data from specific regions of the full domain, as described in section 2.3. Our aims are to both successfully train the neural networks and to study how they generalize to previously unseen regions.

We study the spatiotemporal variability of  $S_x$  and  $\tilde{S}_x$ , by examining snapshots, the time mean, and the standard deviation, shown in Figure 2. Diagnostics are calculated over the full  $512 \times 512$  domain, using the final year of withheld data. Both the spatial and temporal variability of the true  $S_x$  are dominated by the jet dynamics (Figures 2a, 2e, and 2i). In particular, strong meanders which extend eastward from the western boundary are visible. The amplitude of the spatiotemporal variability of  $S_x$  ( $1.4 \times 10^{-6} \text{ m/s}^2$ ) is of similar magnitude to the time mean ( $1.5 \times 10^{-6} \text{ m/s}^2$ ).

All neural networks trained on three different regions, shown in Figure 1a and described in section 2.3, successfully reproduce the spatial patterns of the true  $S_x$ , as shown by snapshots of the predictions  $\tilde{S}_x$  (Figures 2b–2d). Their magnitudes however vary significantly. The predictions of  $f_x(\tilde{\psi}, \mathbf{w}_1)$ , trained on data from the western boundary, are almost identical to the true  $S_x$  and successfully reproduces the correct amplitude and variability (Figures 2b, 2f, and 2j). The neural network  $f_x(\tilde{\psi}, \mathbf{w}_2)$ , trained on data from the eastern boundary, underestimates the magnitude of the true  $S_x$  by approximately 50%, despite reproducing the correct spatial patterns. The predictions of  $f_x(\tilde{\psi}, \mathbf{w}_3)$ , trained on the southern gyre, underestimates the true  $S_x$  by an order of magnitude (Figures 2d, 2h, and 2l).

As the variability of  $S_x$  is dominated by the jet, it is difficult to assess the accuracy of the neural network predictions  $\tilde{S}_x$  in quiescent regions such as the eastern boundary or within the gyres. We therefore calculate the Pearson correlation, a dimensionless quantity, between the true  $S_x$  and the predictions  $\tilde{S}_x$ . The predictions of  $f_x(\tilde{\psi}, \mathbf{w}_1)$  and  $f_x(\tilde{\psi}, \mathbf{w}_2)$  are highly correlated with the truth ( $r > 0.9$ ) within the jet but tend toward zero or negative correlation near the eastern boundary (Figures 2m and 2n). The predictions of  $f_x(\tilde{\psi}, \mathbf{w}_3)$  have a more consistent positive correlation across the gyres and other more quiescent regions, (Figure 2o).





**Figure 2.** Examining the nonlocal prediction ability. Comparisons of the true zonal component of the subfilter momentum forcing  $S_x$ , with the neural networks trained using data from three different regions. The first three rows compare (a–d) snapshots, (e–h) time means, and (i–l) the standard deviation, respectively, while the bottom row (m–o) shows the correlation between the true  $S_x$  and the predictions  $\tilde{S}_x$ . The first column contains the diagnostics using the true zonal subfilter momentum forcing  $S_x$ , while columns two, three, and four use predictions  $\tilde{S}_x$  from the neural networks  $f_x(\bar{\psi}, \mathbf{w}_1)$ ,  $f_x(\bar{\psi}, \mathbf{w}_2)$ , and  $f_x(\bar{\psi}, \mathbf{w}_3)$ , respectively. All diagnostics were produced using the validation data.

We observe similar results for the spatial and temporal variability of  $S_y$ , shown in Figure 3: the variability within the jet dominates, with an amplitude ( $1 \times 10^{-6} \text{ m/s}^2$ ) similar to  $S_x$ . The meandering of the jet again produces complex spatial patterns in  $S_y$ , which when averaged in time, produce a distinct sign change moving across the jet latitudinally. For the predictions  $\tilde{S}_y$ , the neural network trained on the western boundary,  $f_y(\tilde{\psi}, \mathbf{w}_1)$ , most effectively reproduces the true  $S_y$ . However, the time mean of  $f_y(\tilde{\psi}, \mathbf{w}_1)$  (Figure 3f) has a positive bias everywhere in the domain, whereas the time means of  $f_y(\tilde{\psi}, \mathbf{w}_2)$  and  $f_y(\tilde{\psi}, \mathbf{w}_3)$  (Figures 3g and 3h, respectively) do not.

The correlations between  $S_y$  and  $\tilde{S}_y$  are similar to the zonal component:  $f_y(\tilde{\psi}, \mathbf{w}_1)$  and  $f_y(\tilde{\psi}, \mathbf{w}_2)$  are highly correlated ( $r > 0.8$ ) within the jet but not in the gyres. In contrast,  $f_y(\tilde{\psi}, \mathbf{w}_3)$  has a consistently positive correlation across the full domain, despite failing to reproduce the amplitude within the jet. In fact, the correlation of  $f_y(\tilde{\psi}, \mathbf{w}_3)$  within the jet (Figure 3o) is negative ( $r \approx -0.3$ ). The negative correlation implies that the dynamical processes occurring within region 3, the southern gyre, have an opposite effect to the eddy momentum forcing occurring within region 1. The opposing effects of eddies could be an example of regional variation in eddy forcing, as in Waterman and Jayne (2010), who found that whether eddies were driving the large-scale flow or not, depended critically on along-stream position.

Across all neural networks, the correlation decreases at the eastern boundary, which is partly caused by the subfilter momentum forcing being orders of magnitude lower than elsewhere in the domain. The low magnitude of  $S_x$  and  $S_y$  is due to the wave-like behavior of the flow having a larger spatial scale. The larger spatial scale at the eastern boundary leads to little variability at small scales, reducing the eddy momentum forcing to almost zero and therefore causing the performance of neural networks to deteriorate.

Overall, we see that training neural networks on the western boundary is most successful when generalizing to other areas of the domain (in terms of correlations and reproducing the variability). Training on the eastern boundary produced good correlations in the western boundary, but underestimated the magnitude of the eddy forcing by approximately 50%. Training on the southern gyre did not correlate well within the western boundary and underestimated the truth by an order of magnitude.

Hence, to successfully reproduce the correct amplitude and variability across the domain, the training data must contain a diverse range of scale interactions, which here corresponds to training on the most turbulent region. However, training on the turbulent regions can lead to significant net biases in the predictions, as seen in Figure 3f. How to correct for such biases will be discussed in section 4.

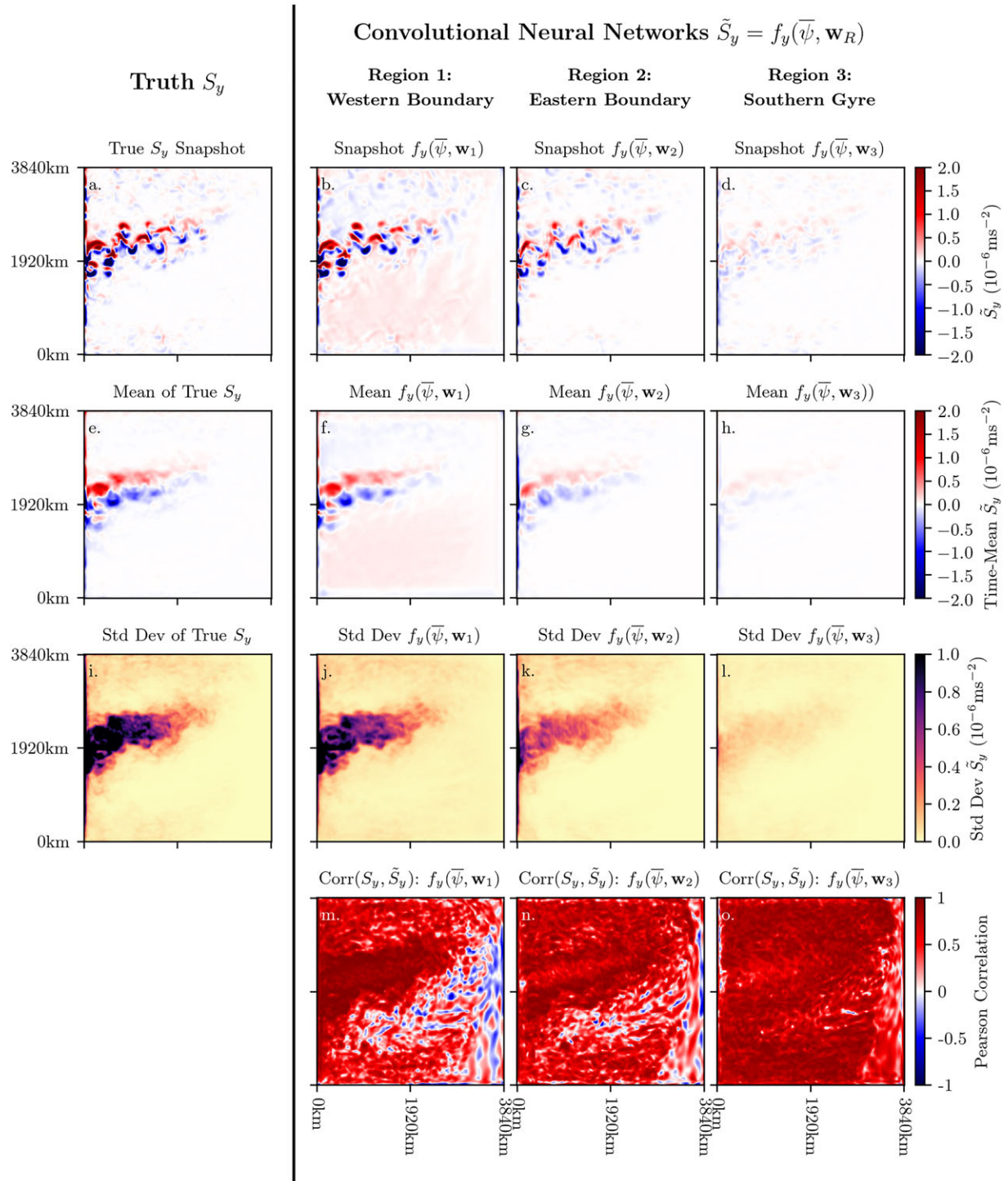
### 3.2. Generalizing to Different Reynolds Numbers

In section 3.1, we investigated how neural networks trained on different regions of the domain generalize to other previously unseen regions. We now test how the neural networks generalize to different regimes, in particular, different Reynolds number. In section 3.1, we found that the neural networks trained on region 1, the western boundary, successfully generalized to different regions; we therefore apply  $f_x(\tilde{\psi}, \mathbf{w}_1)$  to new model data with different wind stress amplitudes and viscosity coefficients to test its performance. We use models with higher and lower wind forcings, to test regimes which are both more and less turbulent than the original model, which had a wind stress amplitude of  $\tau_0 = 0.8 \text{ N/m}^2$  and viscosity  $\nu = 75 \text{ m}^2/\text{s}^2$ .

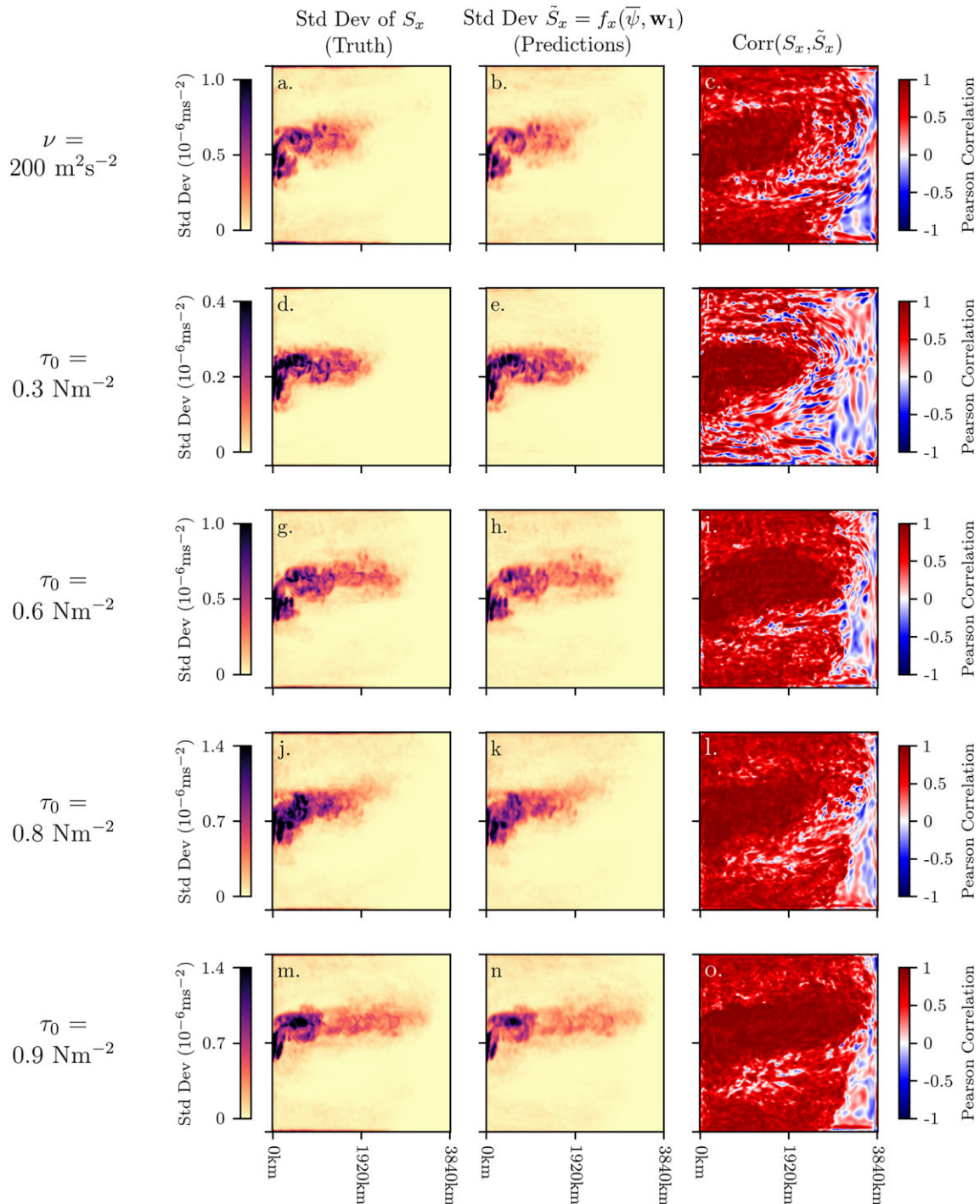
We use the low-pass filter on the upper-layer stream function from each different model run, with the following:  $\nu = 200 \text{ m}^2/\text{s}^2$  and  $\tau_0 = 0.3, 0.6,$  and  $0.9 \text{ N/m}^2$ , and then apply the already trained neural network  $f_x(\tilde{\psi}, \mathbf{w}_1)$  to generate predictions  $\tilde{S}_x$ . The standard deviation of the true  $S_x$ , the standard deviation of the  $f_x(\tilde{\psi}, \mathbf{w}_1)$  predictions  $\tilde{S}_x$ , and the correlation between them are shown in Figure 4.

The neural network  $f_x(\tilde{\psi}, \mathbf{w}_1)$  reproduces the variability within the jet almost exactly, across all runs, as can be seen by comparing the standard deviations in the first and second columns, which represent the standard deviation of the true  $S_x$  and predicted  $\tilde{S}_x$ , respectively. The correlation within the jet remains high ( $r > 0.9$ ) in all runs, including the model with an increased wind forcing ( $\tau_0 = 0.9 \text{ N/m}^2$ ) in Figure 4o. The correlations weaken at the eastern boundary for the lowest wind forcing ( $\tau_0 = 0.3 \text{ N/m}^2$ ), shown in Figure 4f; this may be caused by an increase in the wave-like behavior at the eastern boundary, which is not well captured by the neural networks. In general, the higher the Reynolds number, the better the correlations, that is, more dark red areas of  $r > 0.8$ .

The mean biases of the predictions of the new models are similar in magnitude to the biases of the original model configuration. These biases showed no relationship with the Reynolds number and are therefore not discussed further.



**Figure 3.** The same diagnostics as Figure 2, but for the meridional component of the subfilter momentum forcing: the true  $S_y$  and the predictions  $\tilde{S}_y$  from the neural networks  $f_y(\bar{\psi}, \mathbf{w}_1)$ ,  $f_y(\bar{\psi}, \mathbf{w}_2)$ , and  $f_y(\bar{\psi}, \mathbf{w}_3)$ .



**Figure 4.** Examining the ability to generalize to new regimes: using the trained neural network  $f_x(\bar{\psi}, \mathbf{w}_1)$ , we make predictions for model runs of different viscosities and wind forcings. From each model run, we use 1 year of the upper-layer filtered-stream function to generate predictions  $\tilde{S}_x$  from  $f_x(\bar{\psi}, \mathbf{w}_1)$  to see how they compare to the true  $S_x$ . We study a run of higher viscosity (a–c)  $\nu = 200 \text{ m}^2/\text{s}^2$ , and runs with wind stress amplitude (d–f)  $\tau_0 = 0.3$ , (g–h)  $0.6$ , (g–i)  $0.8$ , and (j–l)  $0.9 \text{ N/m}^2$ . Note that  $f_x(\bar{\psi}, \mathbf{w}_1)$  was trained on a run with  $\nu = 75 \text{ m}^2/\text{s}^2$  and  $\tau_0 = 0.8 \text{ N/m}^2$ , the standard deviation and correlation maps of which are included again here in panels (j), (k), and (l).

### 3.3. Sensitivity of Neural Networks to Undersampling

We have so far trained the neural networks with densely sampled data; that is, we have data at each grid point for both the input and output variables. However, most observational data sets are spatially sparse, for example, Argo floats (Roemmich et al., 2009). We therefore explore the impact of undersampling with a new collection of neural networks trained on region 1 to predict  $S_x$ , but with the training data subsampled. At each time slice of the training data, we randomly sample (without replacement)  $N$  points of the  $40 \times 40$  input variables,  $\tilde{\psi}$ , and output variables  $S_x$ . Using these  $N$  randomly sampled values, we use a cubic interpolation to reconstruct the full  $40 \times 40$  grid point input and output (with a nearest neighbor interpolation for grid points that fall outside the convex hull of the cubic interpolation).

These reconstructed time slices from subsampled data are used to train a new set of neural networks. We vary the number of points  $N$  subsampled from  $>90\%$  to  $<5\%$  of the original 1,600 points of the input and output variables. We have a neural network for each value of  $N$ , the subsampling rate. Using the neural networks trained on undersampled data, we calculate the root-mean-square error (RMSE) on the final year of validation data over the entire domain. The validation data are not subsampled, providing a stronger and more accurate test of the neural network's performance.

The RMSE is shown as a function of percentage of points sampled (Figure 5c). We find that the RMSE increases significantly only when the percentage of spatial points sampled drops below 10% (the error doubles at a subsampling rate of 4.7%). Note that the RMSE is not a monotonic function of percentage of points sampled due to the stochastic nature of the training procedure and the use of a nonlinear interpolation. The spatial map of RMSE of the neural network trained with 18.75% subsampled data (Figure 5b) shows minimal changes relative to the neural network trained on the original (unaltered) training data (Figure 5a). The result further suggests that the use of sparse interpolated observations can be successfully used to accurately train and predict the eddy momentum forcing as shown in sections 3.1 and 3.2.

We also tested an alternative method of undersampling, where the  $40 \times 40$  input and output grid of the neural network is spaced out over the entire domain. In other words, we subsample the input and output variables of the original  $512 \times 512$  grid to a regularly spaced  $40 \times 40$  grid. However, training a convolutional neural network with this methodology did not work and led to severe overfitting (i.e., increasing validation loss during training). The neural networks presented in section 2 learn to take first- and second-order derivatives of the input stream function (see GitHub repository), which correspond to the velocities and velocity shears. Both velocities and velocity shears are important features to provide for accurate predictions of the eddy momentum forcing. By severely subsampling the input stream function, the local information relevant to estimate velocities and velocity shears is lost.

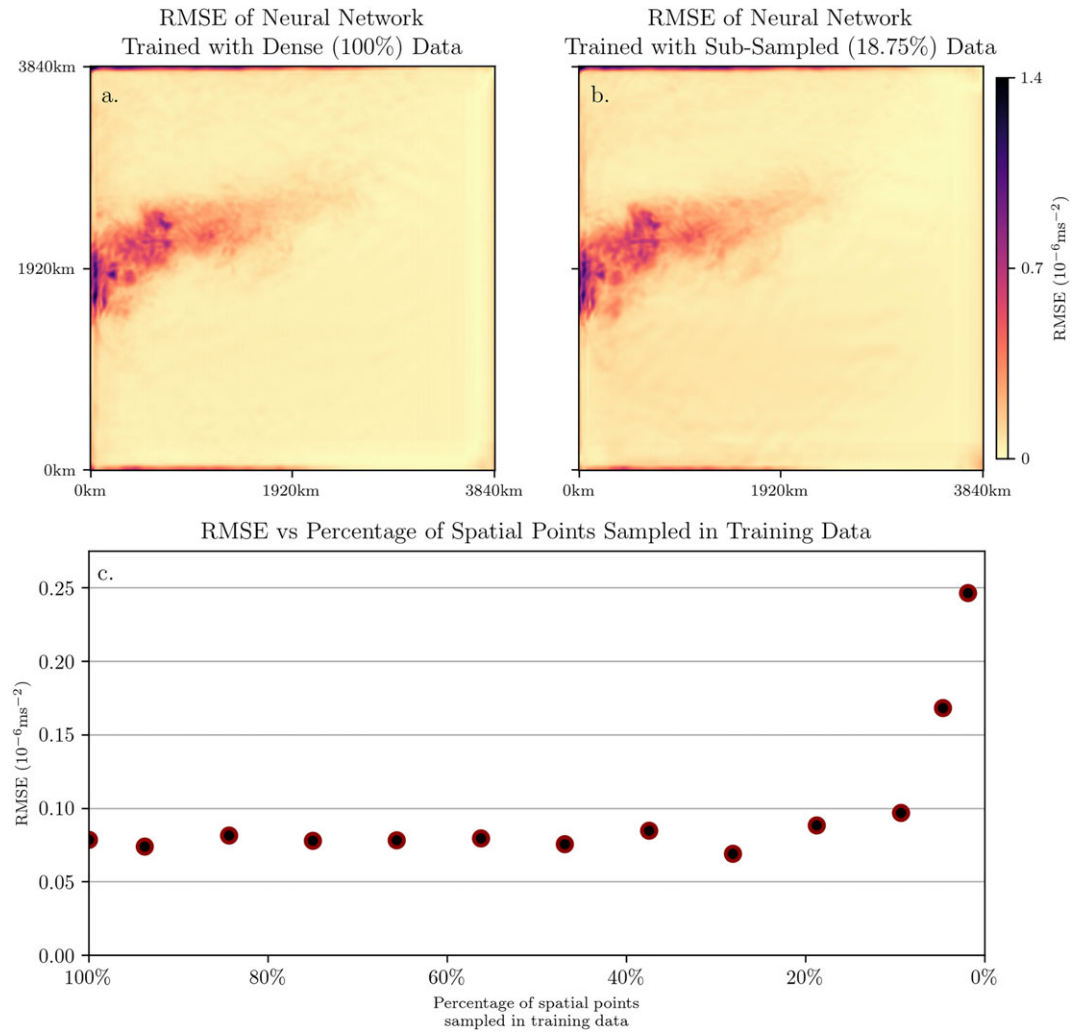
Therefore, the success of training convolutional neural networks on observational data sets will likely depend on the horizontal sampling rate of the product being considered. Our results suggest that high horizontal resolution sampling is necessary to capture any small-scale gradients and avoid the sharp rise in error in Figure 5. We anticipate that the use of (interpolated) data sets with horizontal resolution of  $1^\circ$  or less (e.g., Argo float, altimetry) is needed to train the neural networks and predict the eddy momentum forcing.

## 4. Physically Constrained Neural Networks

We proceed to examine the net input of momentum from the neural network predictions  $\tilde{S}_x$  and  $\tilde{S}_y$ , which should vanish. If neural networks are used to leverage the use of observational data sets and coarse-resolution models, then spurious sources of momentum would violate physical conservation laws. We therefore need to constrain the neural networks to respect the physical properties of the system. Here we diagnose the momentum biases of the neural networks  $f_i(\tilde{\psi}, \mathbf{w}_R)$ , and then explore different methods of imposing conservation of momentum globally.

### 4.1. Momentum Biases

Each subregion (including those used to train the neural networks) may have a nonzero spatially integrated momentum tendency. However, globally, the true subfilter momentum forcing  $\mathbf{S}$  should redistribute momentum and not act as a source or sink, that is,  $\iint \mathbf{S} dx dy = 0$ . We therefore need the neural networks to not introduce spurious sources of momentum, to respect the physical properties of the system. By training each neural network on a subregion, we expect to have imperfect momentum conservation, which will depend upon the particular dynamical processes within each region. For example, if eddies within a par-



**Figure 5.** Determining how undersampling of the training data impacts neural network error. Panel (a) shows the root-mean-square error (RMSE) of the neural network  $f_x(\bar{\psi}, \mathbf{w}_1)$  trained with dense (unaltered) training data, while panel (b) shows the RMSE of the neural network trained with subsampled (18.75%) data. Panel (c) shows the RMSE as a function of the percentage of spatial points sampled at each time slice of the training data. Note that the RMSE is calculated over the full domain during the validation period (the final year of data).

ticular region are driving the mean flow, then we would expect a positive source of momentum locally—a neural network trained on such a region would likely generalize the (local) input of momentum to the rest of the domain. A net source or sink of momentum will manifest as a nonzero bias after spatial averaging.

At a single point in space, the time series of the predictions  $\tilde{S}_x$  and  $\tilde{S}_y$  show that the neural networks trained on regions 1 and 2 track the true  $S_x$  and  $S_y$  closely (Figures 6a and 6b), reproducing a significant proportion (>80%) of the variance. However, if at each time step we spatially average the neural network predictions  $\tilde{S}_x$  and  $\tilde{S}_y$  (Figures 6c and 6d, respectively) over the full domain, we observe significant nonzero biases.

Consider the zonal component of the eddy momentum forcing in Figure 6c:  $f_x(\bar{\psi}, \mathbf{w}_1)$  has a net positive bias, implying a global positive increase of zonal momentum at all times, while both  $f_x(\bar{\psi}, \mathbf{w}_2)$  and  $f_x(\bar{\psi}, \mathbf{w}_3)$  have negative biases, indicating a net decrease in zonal momentum. We can estimate the magnitude of the resulting change in zonal velocity from these net biases, over a period of a year, by assuming  $\Delta u = \langle \tilde{S}_x \rangle \Delta t$ , where  $\langle \rangle$  denotes the spatial average over the full domain. For  $f_x(\bar{\psi}, \mathbf{w}_1)$ ,  $f_x(\bar{\psi}, \mathbf{w}_2)$ , and  $f_x(\bar{\psi}, \mathbf{w}_3)$ , we obtain values of  $\langle \tilde{S}_x \rangle = 0.03, 0.02,$  and  $0.0008$  ( $10^{-6} \text{ m/s}^2$ ), respectively; this leads to zonal velocity changes of  $\Delta u = 0.95, 0.63,$  and  $0.025$  (m/s). These changes are of similar magnitude to the time mean zonal flow, which peaks at approximately 0.9 m/s within the jet core.

There are also significant biases in the predictions of the meridional component  $\tilde{S}_y$ , shown in Figure 6d. The positive bias of  $f_y(\tilde{\psi}, \mathbf{w}_1)$  is visible in the time mean  $\tilde{S}_y$ , shown in Figure 3f. We can again estimate the change in meridional velocities by assuming  $\Delta v = \langle \tilde{S}_y \rangle > \Delta t$ . Using values of  $\langle \tilde{S}_y \rangle = 0.02, -0.01, \text{ and } 0.002$  ( $10^{-6}$  m/s<sup>2</sup>) for  $f_y(\tilde{\psi}, \mathbf{w}_1), f_y(\tilde{\psi}, \mathbf{w}_2), \text{ and } f_y(\tilde{\psi}, \mathbf{w}_3)$ , respectively, leads to the following changes:  $\Delta v = 0.63, -0.31, \text{ and } 0.06$  (m/s). Some of these changes are the same magnitude as the time mean meridional flow.

#### 4.2. Toward Momentum-Conserving Neural Networks

The predictions of neural networks  $f_x(\tilde{\psi}, \mathbf{w}_1)$  and  $f_y(\tilde{\psi}, \mathbf{w}_1)$ , described in section 3.1, correctly reproduce the correct amplitude and variability of the true eddy momentum forcing  $S_x$  and  $S_y$ , as seen in Figures 2 and 3. However, training on region 1 also produced some of the largest nonzero biases in  $\tilde{S}_x$  and  $\tilde{S}_y$  after spatial averaging at each time step. We therefore test whether we can reduce the biases when training on region 1, while preserving the accuracy of predictions from the neural network. We try three approaches (A, B, and C) to reduce the biases identified in Figures 6c and 6d.

- (A) Architecture alteration: Train neural networks on region 1, but with the final fully connected layer modified such that the spatial mean is removed from the final output. The neural networks will therefore be trained to reproduce the subfilter momentum forcing, but with momentum conservation intrinsically embedded, that is, same training data, but altered architecture. The motivation behind this approach is that if the local source of momentum within the  $40 \times 40$  output grid is zero, then this may reduce the global net source of momentum.
- (B) Preprocessing of input: Train on region 1 with the original architecture described in Table 1 but with the spatial mean removed from each  $40 \times 40$  training snapshot of  $S_x$ . In other words, remove the local bias from each training snapshot. If the local source of momentum of each  $40 \times 40$  sample is zero within the training data, then training on such data may reduce local biases when making future predictions on unseen data. However, this does not guarantee that subsequent predictions will have zero local bias. This approach tries to increase the probability of *learning* local momentum conservation.
- (C) Postprocessing of output: Train on region 1 and enforce global momentum conservation after the predictions have been made, that is, no changes to training data or architecture, but with additional processing of the full-domain predictions  $\tilde{S}_x$  and  $\tilde{S}_y$ .

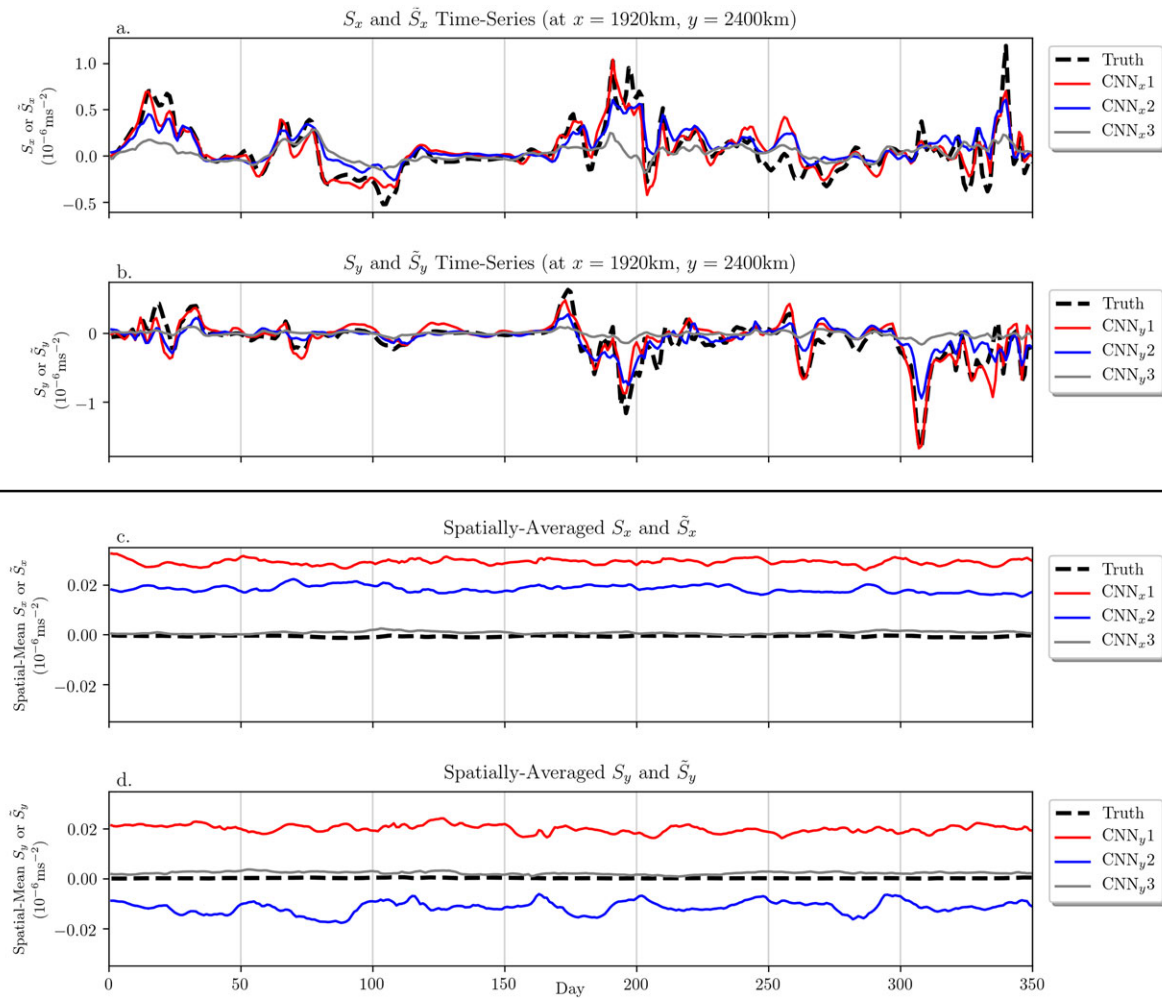
The associated neural networks of each approach are labeled as  $f_i(\tilde{\psi}, \mathbf{w}_1^A), f_i(\tilde{\psi}, \mathbf{w}_1^B), \text{ and } f_i(\tilde{\psi}, \mathbf{w}_1^C)$ , respectively, where  $i = (x, y)$  denotes either the zonal  $S_x$  or meridional  $S_y$  component being predicted.

All neural networks are optimized using the same training parameters given in Table 1. Approach A, which alters the architecture, and approach B, which alters the training data, are enforcing momentum conservation not just globally, but within the  $40 \times 40$  subregion being predicted. This local conservation is useful for enforcing global conservation. However, local conservation may not be desirable if there is convergence of eddy momentum fluxes in a particular region, which can impact the large-scale flow; for example, if eddies are fluxing momentum into the jet at a particular along-stream position, enforcing local conservation in a neural network may lead to missing these effects. Therefore, caution must be taken with restricting architectures in this way.

We now explore the performance of the newly constrained neural networks and the net momentum input relative to that of the original neural networks trained on region 1:  $f_x(\tilde{\psi}, \mathbf{w}_1)$  and  $f_y(\tilde{\psi}, \mathbf{w}_1)$ . The spatial averages of neural networks based on approaches A, B, and C are shown in Figure 7, with the same scale axes as in Figure 6.

Approach B has significant biases of approximately  $-0.01$  and  $-0.015$  ( $10^{-6}$  m/s<sup>2</sup>) in the zonal and meridional components, respectively; the optimization procedure aims to reproduce the *variability* in the training data, and not spatial means; therefore, preprocessing the training data does not remove the biases. Compared to the original neural networks trained on region 1, the biases of approaches A and C are 3 to 5 orders of magnitude lower, in both the zonal and meridional components. The postprocessing approach is exactly zero by construction, while the altered architecture approach A is not exactly zero due to the overlapping tiling procedure. The biases of  $f_x(\tilde{\psi}, \mathbf{w}_1^A)$  and  $f_y(\tilde{\psi}, \mathbf{w}_1^A)$  are approximately  $-0.002$  and  $-0.0005$  ( $10^{-6}$  m/s<sup>2</sup>) which, over the course of a year, would lead to velocity changes of  $\Delta u = -0.06$  and  $\Delta v = -0.01$  (m/s), respectively—now an order of magnitude smaller than the time mean flow.

The correlation maps of all momentum-conserving approaches (not shown) change little from the original correlation maps of  $f_x(\tilde{\psi}, \mathbf{w}_1)$ , and  $f_y(\tilde{\psi}, \mathbf{w}_1)$ , shown in Figures 2m and 3m, respectively. All approaches



**Figure 6.** Panels (a) and (b) show time series of the zonal and meridional components of the subfilter momentum forcing, respectively, at a single point near the middle of the domain. Panels (c) and (d) also show time series of the zonal and meridional components of the subfilter momentum forcing, but this time spatially averaged over the entire domain.

reproduce the correct spatial patterns of the true  $S_y$  and  $\tilde{S}_y$  (e.g., Figure 7 for standard deviations). However, approaches A and B underestimate the amplitude of  $S_x$  and  $S_y$  by approximately 20–30%, whereas there is a little difference between approach C and the truth (<10%).

In summary, approach C of postprocessing successfully enforces momentum conservation, without sacrificing accuracy in the predictions of the eddy momentum forcing. Approach B, altering the training data, was not efficacious at reducing the net biases. The physically constrained architecture of approach A successfully reduced the net bias, but at the expense of 20–30% accuracy. Though further altering of the architecture (e.g., increasing number of convolution layers and filters) or training procedure (decreasing the learning rate, with increased number of training epochs) could reduce this drop in accuracy by countering the restriction placed on the architecture.

### 5. Predicting Subsurface Flow

We have shown that neural networks, by using the filtered-stream function as the input variable, can provide information on unresolved turbulent processes, namely, the subfilter momentum forcing. We have assumed that the filtered-stream function represents some limited set of observations or data from a coarse-resolution ocean model. However, coarse-resolution ocean models still produce data for below the surface, whereas satellite observations do not. Here we address the issue of inferring subsurface information solely from surface fields. Our approach is conceptually similar to Chapman and Charantonis (2017),



which used a form of neural network called a self-organizing map to reconstruct subsurface velocities in the Southern Ocean, using satellite altimetry and Argo float data. Using the QG model data described in section 2.1, we test whether a neural network can predict the middle-layer stream function, using only the surface-filtered-stream function.

We train a new neural network  $\tilde{\psi}_2 = f(\tilde{\psi}_1, \mathbf{W})$  (which has the same architecture as before, but with a different output and weights) to minimize the mean-square-error loss function  $L \propto (\psi_2 - \tilde{\psi}_2)^2$ , where  $\tilde{\psi}_1$  is the filtered-stream function of the upper layer,  $\psi_2$  is the true stream function of the middle layer, and  $\tilde{\psi}_2$  is the neural network predictions. Again, to assess the ability to generalize to unseen regions, we only train the neural network on the western boundary (training region 1). Diagnostics of the true  $\psi_2$  and predictions  $\tilde{\psi}_2$ , including the correlation between them, are shown in Figures 8a–8e. The neural network accurately reproduces the middle-layer time mean and standard deviation of the stream function within the jet region. The neural network accurately reproduces the correct amplitude of the true  $\psi_2$  within the jet, but underestimates the amplitude by  $\approx 50\%$  within the gyres. Independent of the amplitude, the predictions  $\tilde{\psi}_2$  are highly correlated ( $r > 0.8$ ) almost everywhere in the domain with the true  $\psi_2$ .

The decrease in accuracy in the gyres is likely due to only training within the western boundary, where the stream functions of the upper- and middle layers are more tightly coupled due to the strong barotropic nature of the flow. Within the gyres, the barotropic component is not as dominant—this could cause the neural networks to underestimate the amplitude away from the jet. Alternatively, the adjustment time scales of the upper and middle layers are not the same, which perhaps requires more training data in order to capture interactions over longer time scales.

We take the approach one step further, by predicting the bottom-layer stream function, using the same neural network and its weights  $f(\tilde{\psi}_1, \mathbf{W})$ , but now using the predictions of the middle-layer stream function as the input, that is,  $\tilde{\psi}_3 = f(\tilde{\psi}_2, \mathbf{W})$ . We test whether a neural network trained to predict the middle-layer stream function can provide any information on the bottom-layer stream function (without retraining), by inputting the middle-layer stream function as an input. Mathematically, this is written as  $\tilde{\psi}_3 = f(f(\tilde{\psi}_1, \mathbf{W}), \mathbf{W})$ .

Diagnostics of the true ( $\psi_3$ ) and predicted ( $\tilde{\psi}_3$ ) bottom-layer stream function are shown in Figures 8f–8j. Despite a moderate correlation of  $r \approx 0.5$  across the domain, the predictions fail to reproduce the correct time mean, which has a circulation in the opposite direction to the truth. This is due to the neural network being trained to predict the middle-layer flow, which on average is more aligned with the upper layer. Therefore, when the neural network is given the middle-layer stream function as an input, it predicts the bottom-layer flow as on average being in the same direction, which is not the case. The neural network also has not been trained to predict the effects of the additional bottom drag, decreasing the accuracy further—more data could improve this issue, as the longer time scales associated with bottom drag may be absent from the training data set.

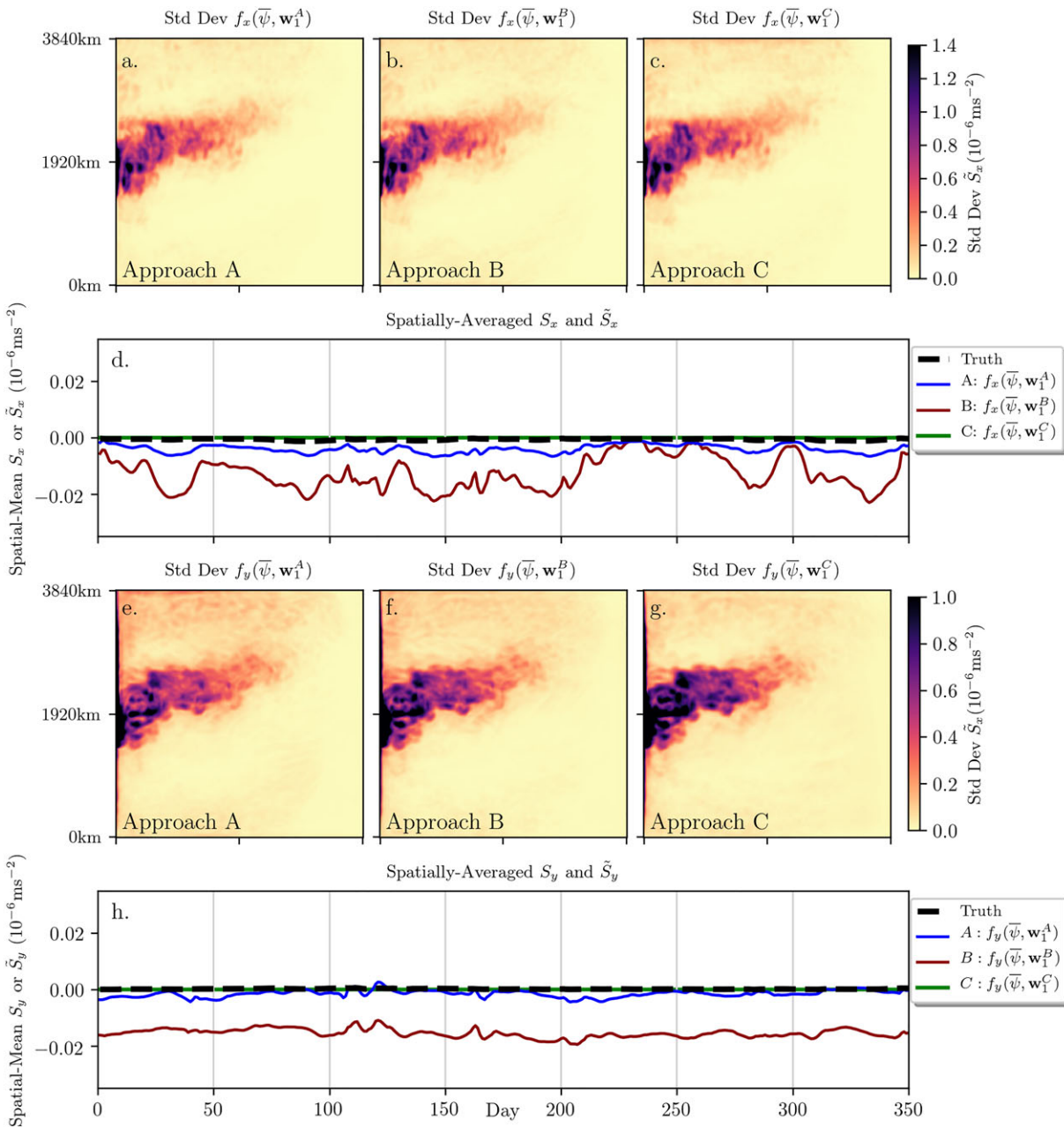
An alternative approach would be to train a new neural network to map directly from the surface flow to the bottom-layer flow, that is,  $\tilde{\psi}_3 = f(\tilde{\psi}_1, \mathbf{W})$ . Having separate neural networks for the middle and bottom layers, you could then reconstruct the flow at all depths using just information at the surface (although an additional neural network does increase computational costs). Independent of the abyssal flow however, we have shown that neural networks can provide information on the flow at intermediate depths.

## 6. Conclusions and Discussion

### 6.1. Summary

In this study, we have demonstrated as a proof of concept that machine learning algorithms can provide information on unresolved turbulent processes, when given a smoothed view of the dynamics (i.e., the filtered-stream function). We degrade data from a high-resolution eddy-resolving QG model using a spatial low-pass filter and train convolutional neural networks to predict the relationship between turbulent processes and their effect on the large-scale flow, that is, the eddy momentum forcing. Our results show that convolutional neural networks can successfully represent both the spatial and temporal variability of the eddy momentum forcing.

We determine how neural networks trained on one area of the domain, perform in other previously unseen areas (Figures 2 and 3), representing when observational data are limited to only particular regions, for example, mooring data (Hogg, 1992) or gliders (Davis et al., 2008; Rudnick et al., 2004). Training on a sub-



**Figure 7.** The standard deviation and spatial-average time series of the predictions  $\tilde{S}_x$  and  $\tilde{S}_y$  of the momentum conserving approaches A, B, and C. Panels (a), (b), and (c) show the standard deviation of  $\tilde{S}_x$  from  $f_x(\bar{\psi}, \mathbf{w}_1^A)$ ,  $f_x(\bar{\psi}, \mathbf{w}_1^B)$ , and  $f_x(\bar{\psi}, \mathbf{w}_1^C)$ , respectively, while panels (e), (f), and (g) show the standard deviation of  $\tilde{S}_y$  from  $f_y(\bar{\psi}, \mathbf{w}_1^A)$ ,  $f_y(\bar{\psi}, \mathbf{w}_1^B)$ , and  $f_y(\bar{\psi}, \mathbf{w}_1^C)$ , respectively. The spatial averages of these predictions  $\tilde{S}_x$  and  $\tilde{S}_y$  are shown in panels (d) and (h).

region tests the sensitivity of the neural network performance to the underlying physical processes. We find that the region on which the neural network is trained significantly impacts the accuracy, as well as the mean bias, which impacts momentum conservation. In particular, training on the least energetically active region, the southern gyre, leads to the lowest accuracy; these neural networks could not reproduce the variability in more energetic regions, such as within the meandering jet. However, training on the western boundary leads to the best generalization, in terms of reproducing the correct amplitude of the eddy momentum forcing in the rest of the domain.

The variation in performance between regions implies that training on the most turbulent region leads to the best performing neural networks for eddy momentum forcing prediction. It is possible that data from

the most turbulent regions exhibits the highest variance or contains a more diverse range of scale interactions. However, two regions may be as turbulent or energetically active as each other, but the nature of the eddy-mean flow interactions within them may differ. For example, Waterman and Jayne (2010) showed that in an idealized model the effect of eddies on the mean flow depended critically on along-stream position: up-stream eddies are generated by an unstable jet, while downstream the eddies drive the time mean circulation. Therefore, training neural networks on different along-stream positions may lead to different dynamical processes being learned, despite both regions being energetically active. Here we have shown how the performance varies between regions of differing energetic activity, but how the specific effects of eddies—for example, driving the mean flow, versus eddies extracting momentum and energy from the jet—impacts the neural network performance remains to be determined.

Without further training, we show that a neural network trained on one QG model configuration generalizes exceedingly well to QG models with different viscosity coefficients and wind forcings (Figure 4). The neural network within the jet reproduces the correct spatiotemporal variability (<10% error) in all configurations, and the correlation between the predicted  $\tilde{S}_x$  and the true  $S_x$  within the gyres increases with the Reynolds number of the model configuration. While the neural networks do not conserve momentum globally (Figure 6c and 6d), we show that momentum conservation can be enforced without a significant reduction in accuracy (Figure 7), through either a physically constrained architecture or postprocessing of the predictions.

We also show that a new neural network can be trained to predict the middle-layer stream function, using only the upper-layer stream function as the input, that is, predicting the flow at depth using information at the surface (Figure 8). The highest accuracy occurs where the barotropic component of the flow is most dominant, which coincides with a strong zonal mean flow. However, when using the stream function to predict the bottom-layer stream function, the neural network captures some of the variability, but fails to replicate the time mean of the true bottom-layer stream function  $\psi_3$  (Figure 8), primarily due to the presence of bottom drag.

### 6.2. Implications for Leveraging Observations

Our work has implications for inference from sparse observations. While previous studies have used machine learning to leverage observational data sets (Chapman & Charantonis, 2017; Giglio et al., 2018; Su et al., 2018), the present work demonstrates that convolutional neural networks in particular are an excellent tool for such tasks. Neural networks should be further tested and exploited in the future for data inference due to

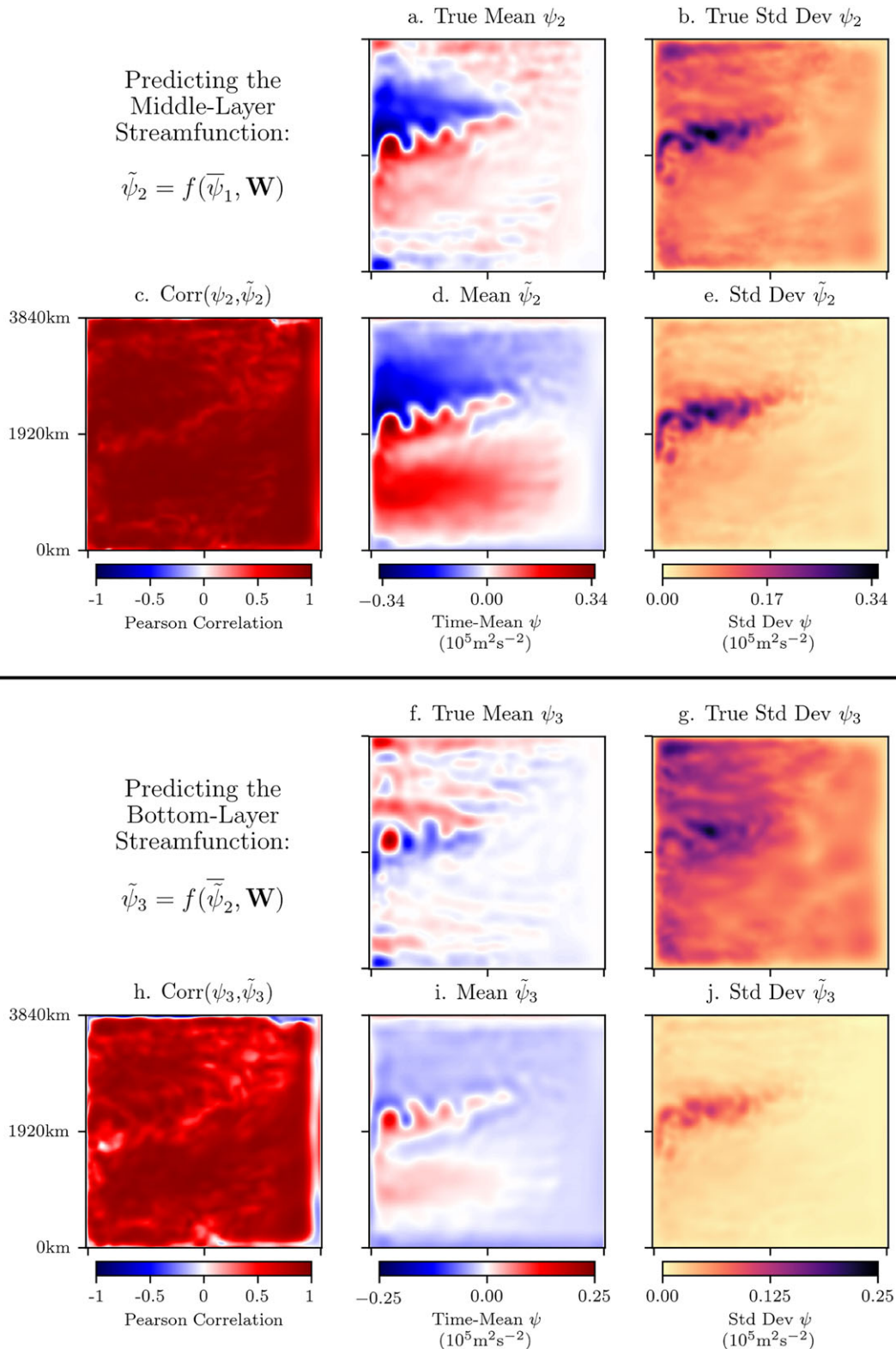
- their resilience, such that accurate predictions for the full domain can be generated by training on a subregion;
- their generalization to different external forcings, without any further training such that predictions outside the regime trained on can be successful; and
- their ability to be successfully trained with undersampled data (Figure 5).

Collectively, these results suggest that sparse interpolated observational data sets can be leveraged by such data-driven techniques. For example, satellite altimetry data can be used to predict the subsurface flow; or data from moorings deployed in Drake Passage as part of the Diapycnal and Isopycnal Mixing Experiment in the Southern Ocean can be used to infer eddy momentum or heat flux divergences in other parts of the Southern Ocean. In addition, data sets from Argo floats (Chapman & Sallée, 2017), mooring data, ADCPs, and SSH from altimetry, could be combined to reconstruct physically and biogeochemically important quantities such as energy reservoirs, or air-sea fluxes, interior transport and/or storage of heat, and carbon and oxygen in the ocean (Giglio et al., 2018; Su et al., 2018).

### 6.3. Implications for Parameterizations

Although we have motivated our study through the leverage of observations and coarse-resolution model data, our results have implications for eddy parameterizations of momentum and more generally for subgrid parameterizations. As discussed previously, machine learning has been used to parameterize unresolved processes in the atmosphere (Brenowitz & Bretherton, 2018; Gentine et al., 2018; Jiang et al., 2018; O'Gorman & Dwyer, 2018).

We have shown that neural networks can successfully represent the spatiotemporal variability of the eddy momentum forcing, implying potential for data-driven oceanic turbulence closures in the future, as sug-



**Figure 8.** Predicting the middle- and bottom-layer stream functions  $\psi_2$  and  $\psi_3$  using the upper-layer filtered-stream function  $\bar{\psi}_1$ . We first train a new neural network to predict  $\psi_2$  from  $\bar{\psi}_1$ , that is,  $\psi_2 = f(\bar{\psi}_1, \mathbf{W})$ ; diagnostics of the true  $\psi_2$  and the predictions  $\tilde{\psi}_2$  are shown in the top-half of the figure (a–e). We then take the same neural network that was trained to predict  $\psi_2$  from  $\bar{\psi}_1$ , and now predict the bottom layer stream function  $\psi_3$  using the predicted middle-layer stream function as the input, that is,  $\psi_3 = f(\tilde{\psi}_2, \mathbf{W})$ ; the diagnostics of the true  $\psi_3$  and the predictions  $\tilde{\psi}_3$  are shown in the bottom-half of the figure (f–j).

gested by Zanna et al. (2018). The generalization ability of the neural networks shows that only a limited amount of observations or high-resolution model data may be needed to successfully represent subgrid-scale processes. While the CNNs are successful at representing relationship between the eddy momentum forcing and their effect on the resolved flow, the low-resolution climate models might have biases that are too severe (e.g., weak transport and velocity shears) to lead to a successful representation of the eddy momentum forcing from CNNs as trained here. Nonetheless, our results also show that they perform very well for different amplitudes of forcing and dissipation. Therefore, until the CNNs are implemented into a coarse-resolution ocean model, their success in improving numerical simulations is speculative but deserves to be investigated.

Whether neural networks are being used to leverage observations, or more importantly to construct a data-driven eddy parameterization, caution must be taken to ensure that the laws of physics are respected. More work into physically constrained machine learning algorithms is crucial, and successful applications of data-driven techniques should incorporate physical knowledge. Indeed, the neural network turbulence model of Ling, Kurzawski, and Templeton 2016 out-performed more simple linear models only when Galilean invariant stress tensors from Pope (1975) were used, which are also a key ingredient of the eddy parameterization proposed by Anstey and Zanna (2017). As previously discussed, we successfully enforce global momentum conservation in the present work, such that future implementations of data-driven parameterizations, despite being semiempirical, can be altered to respect physical principles. Specifically, physical constraints can be incorporated into the architecture of the predictive algorithms.

One potential disadvantage of convolutional neural networks may be the computational cost of the matrix operations of each convolution layer to make a prediction given an input. The total time complexity (ignoring any fully connected layers) of a CNN (He & Sun, 2015) is given by  $\mathcal{O}(\sum_l^d n_{l-1} \cdot s_l^2 \cdot n_l \cdot m_l^2)$ , where  $d$  is the total number of convolution layers,  $l$  is the index of a convolution layer,  $n_l$  is the number of filters,  $s_l$  is the filter size, and  $m_l$  is the size of the output feature map. The time complexity is larger than that of a traditional eddy closure (e.g., a simple Laplacian dissipation of momentum which only involves a few matrix additions and subtractions). One way to reduce the time complexity is to instead use depth-wise separable convolution layers (e.g., Howard et al., 2017), which treat the input channels of a convolution layer more independently. This reduces the number of parameters and hence computational cost. An alternative way of reducing time complexity is to simply reduce the sizes of the input and outputs; that is, make predictions for a region smaller than  $40 \times 40$  grid points. The amount of information available to make predictions is therefore reduced. The computational cost is an area which needs addressing if CNNs are to be routinely implemented in models in the future. However, unlike other parameterizations, the training of the neural networks is only done once.

#### 6.4. Future Work

Our study is a step toward using convolutional neural networks to extend the reach of currently available observational or model data. Our proof-of-concept study was conducted in an idealized QG model. The next stage involves training neural networks on actual observational data sets (as described in section 6.2) or on more realistic model data (e.g., a  $1/40^{\text{th}}$ ° global model which resolves the mesoscale and submesoscale eddy fields, such as in Rocha et al., 2016).

We used 9 years of data to train the neural networks and 1 year for validation. Gentine et al. (2018) showed, with regard to parameterizing convection with neural networks, that the training data set could be reduced in size from 12 to 3 months, with little change in the overall mean-square error. The sensitivity our neural networks to reductions in the amount of training data needs to be systematically explored. We have only determined the impact of spatial undersampling on the neural networks. Further work is needed to determine the impact of using a few number of time slices (e.g., using 3 years of training data as opposed to 9 years).

Training on the western boundary produces the best performance. However, the high skill within the jet does not fully translate to high skill in all parts of the gyres. The best correlations in the gyres occur instead when training on the southern gyre, and not the western or eastern boundaries (Figure 2 and 3). This implies that there may be an optimal combination of the predictions of the neural networks trained on different regions, in order to produce the best overall generalization and potentially include nonlocal effects. For example, each neural network has a weight  $a_i$ , and the optimal predictions for the full domain is a combination of all

neural networks

$$\tilde{S}_x^{\text{OPT}} = \sum_i^N a_i f_x(\tilde{\psi}, \mathbf{w}_i), \quad (10)$$

where the summation is over all regions, and  $\tilde{S}_x^{\text{OPT}}$  is the corresponding optimal prediction (with an analogous  $\tilde{S}_y^{\text{OPT}}$  for the meridional component). Combining predictions from multiple neural networks in this manner could be a useful way of capturing the distinct eddy-mean flow interactions observed by Waterman and Jayne (2010). Alternatively, if the computational resources are available, you could train a single neural network on data from all three regions, in the hope that it *remembers* the physical processes occurring in each region. The risk with this approach is that one loses specialization, and the skill reduces as the single neural network simply *averages* the effects of the three regions together. We will attempt to implement the neural networks (as trained here, or as a combination of neural networks) into a coarse-resolution version of the QG model to test their performance as a subgrid-scale parametrization.

Although this study is a proof of concept, the merging of data-driven methods with physical knowledge has the potential to change the way the physics of the ocean are studied, including the designs of future parameterizations. The combination of physical theory and machine learning could prove more effective than either component in isolation.

#### Acknowledgments

This study was funded by the Natural Environment Research Council (NERC). We thank PierGianLuca Porta Mana, who conducted the high-resolution PEQUOD model simulations. Thank you to Robert Fraser, Tomos David, and Ryan Abernathey for their helpful discussions during the development of this work and to two anonymous reviewers for their comments that helped improve this manuscript. The trained Keras neural networks, their training histories, and the Python code used to produce this paper can all be found in the following GitHub repository: <https://github.com/TomBolton/DeepEddy>.

#### References

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., et al. (2016). Tensorflow: A system for large-scale machine learning, *OSDI* (Vol. 16, pp. 265–283). Savannah, GA: USENIX Association.
- Abernathey, R. P., & Marshall, J. (2013). Global surface eddy diffusivities derived from satellite altimetry. *Journal of Geophysical Research: Oceans*, *118*, 901–916. <https://doi.org/10.1002/jgrc.20066>
- Anderson, G. J., & Lucas, D. D. (2018). Machine learning predictions of a multiresolution climate model ensemble. *Geophysical Research Letters*, *45*, 4273–4280. <https://doi.org/10.1029/2018GL077049>
- Anstey, J. A., & Zanna, L. (2017). A deformation-based parametrization of ocean mesoscale eddy Reynolds stresses. *Ocean Modelling*, *112*, 99–111.
- Arbic, B. K., Polzin, K. L., Scott, R. B., Richman, J. G., & Shriver, J. F. (2013). On eddy viscosity, energy cascades, and the horizontal resolution of gridded satellite altimeter products. *Journal of Physical Oceanography*, *43*(2), 283–300.
- Berloff, P. S. (2005). On dynamically consistent eddy fluxes. *Dynamics of Atmospheres and Oceans*, *38*(3-4), 123–146.
- Brenowitz, N. D., & Bretherton, C. S. (2018). Prognostic validation of a neural network unified physics parameterization. *Geophysical Research Letters*, *45*, 6289–6298. <https://doi.org/10.1029/2018GL078510>
- Chapman, C., & Charantonis, A. A. (2017). Reconstruction of subsurface velocities from satellite observations using iterative self-organizing maps. *IEEE Geoscience and Remote Sensing Letters*, *14*(5), 617–620.
- Chapman, C., & Sallée, J.-B. (2017). Can we reconstruct mean and eddy fluxes from argo floats? *Ocean Modelling*, *120*, 83–100.
- Chelton, D. B., Schlax, M. G., Samelson, R. M., & de Szoeke, R. A. (2007). Global observations of large oceanic eddies. *Geophysical Research Letters*, *34*, L15606. <https://doi.org/10.1029/2007GL030812>
- Chollet, F. (2015). Keras, <https://keras.io>
- Davis, R. E., Ohman, M. D., Rudnick, D. L., & Sherman, J. T. (2008). Glider surveillance of physics and biology in the southern California current system. *Limnology and Oceanography*, *53*(5part2), 2151–2168.
- Dong, C., Loy, C. C., He, K., & Tang, X. (2016). Image super-resolution using deep convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, *38*(2), 295–307.
- Durand, M., Fu, L.-L., Lettenmaier, D. P., Alsdorf, D. E., Rodriguez, E., & Esteban-Fernandez, D. (2010). The surface water and ocean topography mission: Observing terrestrial surface water and oceanic submesoscale eddies. *Proceedings of the IEEE*, *98*(5), 766–779.
- Esteves, J. T., de Souza Rolim, G., & Ferraudo, A. S. (2018). Rainfall prediction methodology with binary multilayer perceptron neural networks. *Climate Dynamics*, 1–13.
- Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., & Yacalis, G. (2018). Could machine learning break the convection parameterization deadlock? *Geophysical Research Letters*, *45*, 5742–5751. <https://doi.org/10.1029/2018GL078202>
- Giglio, D., Lyubchich, V., & Mazloff, M. (2018). Estimating oxygen in the Southern Ocean using argo temperature and salinity. *Journal of Geophysical Research: Oceans*, *123*, 4280–4297. <https://doi.org/10.1029/2017JC013404>
- Goodfellow, I., Bengio, Y., Courville, A., & Bengio, Y. (2016). *Deep learning* (Vol. 1). Cambridge: MIT Press.
- Greatbatch, R., Zhai, X., Claus, M., Czeschel, L., & Rath, W. (2010). Transport driven by eddy momentum fluxes in the Gulf Stream extension region. *Geophysical Research Letters*, *37*, L24401. <https://doi.org/10.1029/2010GL045473>
- Greatbatch, R. J., Zhai, X., Kohlmann, J.-D., & Czeschel, L. (2010). Ocean eddy momentum fluxes at the latitudes of the Gulf Stream and the Kuroshio extensions as revealed by satellite data. *Ocean Dynamics*, *60*(3), 617–628.
- Hallberg, R. (2013). Using a resolution function to regulate parameterizations of oceanic mesoscale eddy effects. *Ocean Modelling*, *72*, 92–103.
- He, K., & Sun, J. (2015). Convolutional neural networks at constrained time cost. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 5353–5360). IEEE, Boston, MA.
- Hewitt, H. T., Roberts, M. J., Hyder, P., Graham, T., Rae, J., Belcher, S. E., et al. (2016). The impact of resolving the Rossby radius at mid-latitudes in the ocean: Results from a high-resolution version of the MET Office GC2 coupled model. *Geoscientific Model Development*, *9*(10), 3655–3670.
- Hogg, N. G. (1992). On the transport of the Gulf Stream between Cape Hatteras and the Grand Banks. *Deep Sea Research Part A. Oceanographic Research Papers*, *39*(7-8), 1231–1246.

- Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., & Adam, H. (2017). Mobilenets: Efficient convolutional neural networks for mobile vision applications. arXiv preprint arXiv:1704.04861.
- Jiang, G.-Q., Xu, J., & Wei, J. (2018). A deep-learning algorithm of neural network for the parameterization of typhoon-ocean feedback in typhoon forecast models. *Geophysical Research Letters*, *45*, 3706–3716. <https://doi.org/10.1002/2018GL077004>
- Jochum, M., Danabasoglu, G., Holland, M., Kwon, Y.-O., & Large, W. (2008). Ocean viscosity and climate. *Journal of Geophysical Research*, *113*, C06017. <https://doi.org/10.1029/2007JC004515>
- Kang, D., & Curchitser, E. N. (2015). Energetics of eddy-mean flow interactions in the Gulf Stream region. *Journal of Physical Oceanography*, *45*(4), 1103–1120.
- Keating, S. R., Majda, A. J., & Smith, K. S. (2012). New methods for estimating ocean eddy heat transport using satellite altimetry. *Monthly Weather Review*, *140*(5), 1703–1722.
- Keating, S. R., & Smith, K. S. (2015). Upper ocean flow statistics estimated from superresolved sea-surface temperature images. *Journal of Geophysical Research: Oceans*, *120*, 1197–1214. <https://doi.org/10.1002/2014JC010357>
- Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980.
- Kjellsson, J., & Zanna, L. (2017). The impact of horizontal resolution on energy transfers in global ocean models. *Fluids*, *2*(3), 45.
- Klambauer, G., Unterthiner, T., Mayr, A., & Hochreiter, S. (2017). Self-normalizing neural networks. In *Advances in neural information processing systems* (pp. 972–981). NIPS, Long Beach, CA.
- Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems* (pp. 1097–1105). NIPS, Lake Tahoe, NV.
- Kutz, J. N. (2017). Deep learning in fluid dynamics. *Journal of Fluid Mechanics*, *814*, 1–4.
- Le Traon, P.-Y., & Morrow, R. (2001). Ocean currents and eddies. In L.-L. Fu & A. Cazenave (Eds.), *Satellite Altimetry and Earth Sciences: A Handbook for Techniques and Applications* (pp. 171–210). San Diego: Academic Press.
- Le Traon, P., Nadal, F., & Ducet, N. (1998). An improved mapping method of multisatellite altimeter data. *Journal of Atmospheric and Oceanic Technology*, *15*(2), 522–534.
- LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, *521*(7553), 436.
- Ling, J., Jones, R., & Templeton, J. (2016). Machine learning strategies for systems with invariance properties. *Journal of Computational Physics*, *318*, 22–35.
- Ling, J., Kurzawski, A., & Templeton, J. (2016). Reynolds averaged turbulence modelling using deep neural networks with embedded invariance. *Journal of Fluid Mechanics*, *807*, 155–166.
- Mana, P. P., & Zanna, L. (2014). Toward a stochastic parameterization of ocean mesoscale eddies. *Ocean Modelling*, *79*, 1–20.
- McGovern, A., Elmore, K. L., Gagne, D. J., Haupt, S. E., Karstens, C. D., Lagerquist, R., Smith, T., & Williams, J. K. (2017). Using artificial intelligence to improve real-time decision-making for high-impact weather. *Bulletin of the American Meteorological Society*, *98*(10), 2073–2090.
- Moeng, C.-H. (1984). A large-eddy-simulation model for the study of planetary boundary-layer turbulence. *Journal of the Atmospheric Sciences*, *41*(13), 2052–2062.
- Morrow, R., Coleman, R., Church, J., & Chelton, D. (1994). Surface eddy momentum flux and velocity variances in the Southern Ocean from Geosat altimetry. *Journal of Physical Oceanography*, *24*(10), 2050–2071.
- O’Gorman, P. A., & Dwyer, J. G. (2018). Using machine learning to parameterize moist convection: Potential for modeling of climate, climate change and extreme events. *Journal of Advances in Modelling Earth Systems*, *10*, 2548–2563. <https://doi.org/10.1029/2018MS001351>
- Pathak, J., Hunt, B., Girvan, M., Lu, Z., & Ott, E. (2018). Model-free prediction of large spatiotemporally chaotic systems from data: A reservoir computing approach. *Physical Review Letters*, *120*(2), 024102.
- Pathak, J., Wikner, A., Fussell, R., Chandra, S., Hunt, B. R., Girvan, M., & Ott, E. (2018). Hybrid forecasting of chaotic processes: Using machine learning in conjunction with a knowledge-based model. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, *28*(4), 041101.
- Pope, S. (1975). A more general effective-viscosity hypothesis. *Journal of Fluid Mechanics*, *72*(2), 331–340.
- Rocha, C. B., Gille, S. T., Chereskin, T. K., & Menemenlis, D. (2016). Seasonality of submesoscale dynamics in the Kuroshio extension. *Geophysical Research Letters*, *43*, 11–304. <https://doi.org/10.1002/2016gl071349>
- Roemmich, D., Johnson, G. C., Riser, S., Davis, R., Gilson, J., Owens, W. B., Garzoli, S. L., Schmid, C., & Ignaszewski, M. (2009). The Argo program: Observing the global ocean with profiling floats. *Oceanography*, *22*(2), 34–43.
- Rudnick, D. L., Davis, R. E., Eriksen, C. C., Fratantoni, D. M., & Perry, M. J. (2004). Underwater gliders for ocean research. *Marine Technology Society Journal*, *38*(2), 73–84.
- Sagaut, P. (2006). *Large eddy simulation for incompressible flows: An introduction*. Verlag Berlin Heidelberg: Springer Science & Business Media.
- Scott, R. B., & Wang, F. (2005). Direct evidence of an oceanic inverse kinetic energy cascade from satellite altimetry. *Journal of Physical Oceanography*, *35*(9), 1650–1666.
- Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556.
- Su, H., Li, W., & Yan, X.-H. (2018). Retrieving temperature anomaly in the global subsurface and deeper ocean from satellite observations. *Journal of Geophysical Research: Oceans*, *123*, 399–410. <https://doi.org/10.1002/2017JC013631>
- Tracey, B. D., Duraisamy, K., & Alonso, J. J. (2015). A machine learning strategy to assist turbulence model development. In *53rd AIAA Aerospace Sciences Meeting* (pp. 1287). Kissimmee, FL: Stanford University.
- 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*, *474*(2213), 20170844.
- Waterman, S., Hogg, N. G., & Jayne, S. R. (2011). Eddy-mean flow interaction in the Kuroshio extension region. *Journal of Physical Oceanography*, *41*(6), 1182–1208.
- Waterman, S., & Jayne, S. R. (2010). Eddy-mean flow interactions in the along-stream development of a western boundary current jet: An idealized model study. *Journal of Physical Oceanography*, *41*, 682–707.

Zanna, L., Brankart, J. M., Huber, M., Leroux, S., Penduff, T., & Williams, P. D. (2018). Uncertainty and scale interactions in ocean ensembles: From seasonal forecasts to multidecadal climate predictions. *Quarterly Journal of the Royal Meteorological Society*. <https://doi.org/10.1002/qj.3397>

Zanna, L., Mana, P. P., Anstey, J., David, T., & Bolton, T. (2017). Scale-aware deterministic and stochastic parametrizations of eddy-mean flow interaction. *Ocean Modelling*, *111*, 66–80.

### Erratum

The figure captions for Figures 5, 6, 7, and 8 were erroneously transposed as the result of a typesetting error. The error has been corrected, and this may be considered the authoritative version of record.