

1 Applications of Deep Learning to Ocean Data Inference and 2 Sub-Grid Parameterisation

3 **Thomas Bolton¹, Laure Zanna¹**

4 ¹Department of Physics, University of Oxford, UK.

5 **Key Points:**

- 6 We successfully use convolutional neural networks to predict unresolved turbulent
7 processes and sub-surface velocities.
- 8 The neural networks generalise to different: regions, dynamical regimes and forc-
9 ing.
- 10 Global momentum conservation for eddy parametrization can be respected without
11 sacrificing accuracy.

12 **Abstract**

13 Oceanographic observations are limited by sampling rates, while ocean models are lim-
 14 ited by finite resolution and high viscosity and diffusion coefficients. Therefore both data
 15 from observations and ocean models lack information at small- and fast-scales. Methods
 16 are needed to either extract information, extrapolate, or up-scale existing oceanographic
 17 datasets, to account for or represent unresolved physical processes. Here we use machine
 18 learning to leverage observations and model data by predicting unresolved turbulent pro-
 19 cesses and sub-surface flow fields. As a proof-of-concept, we train convolutional neural
 20 networks on degraded-data from a high-resolution quasi-geostrophic ocean model. We
 21 demonstrate that convolutional neural networks successfully replicate the spatio-temporal
 22 variability of the sub-grid eddy momentum forcing, are capable of generalising to a range
 23 of dynamical behaviours, and can be forced to respect global momentum conservation.
 24 The training data of our convolutional neural networks can be sub-sampled to 10-20% of
 25 the original size without a significant decrease in accuracy. We also show that the sub-
 26 surface flow field can be predicted using only information at the surface (e.g., using only
 27 satellite altimetry data). Our study indicates that data-driven approaches can be exploited
 28 to predict both sub-grid and large-scale processes, while respecting physical principles,
 29 even when data is limited to a particular region or external forcing. Our in-depth study
 30 presents evidence for the successful design of ocean eddy parameterisations for implemen-
 31 tation in coarse-resolution climate models.

32 **1 Introduction**

33 Satellite observations have produced a wealth of information on the ocean circula-
 34 tion [Morrow *et al.*, 1994; *Le Traon and Morrow*, 2001; *Scott and Wang*, 2005; *Chelton*
 35 *et al.*, 2007; *Greatbatch et al.*, 2010b; *Abernathy and Marshall*, 2013]. However raw satel-
 36 lite altimetry data sub-samples the ocean, and does not measure sub-surface quantities.
 37 Temporally measurements at the same location are made twice every orbital cycle, while
 38 the spatial sampling depends upon the distance between ground tracks. To improve the
 39 sub-sampling rates, measurements from multiple satellites are combined [*Le Traon et al.*,
 40 1998] to produce an optimal estimate.

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

50 Similar to satellite observations, Ocean General Circulation Models (OGCM) are
 51 useful for studying ocean dynamics. However, high-resolution models are computationally
 52 expensive, and the current resolution of models is not high enough to fully resolve the first
 53 baroclinic deformation radius at mid-latitudes [Hallberg, 2013]. Also, due to their finite
 54 resolution, they require large viscosity and diffusion coefficients in order to remain nu-
 55 merically stable [Jochum *et al.*, 2008]. The combination of finite-resolution and artificially
 56 high viscosity, diffuses momentum and smooths out features such as jets and mesoscale
 57 eddies [Hewitt *et al.*, 2016; Kjellsson and Zanna, 2017].

58 Therefore both observations and models are missing the interactions of oceanic tur-
 59 bulence at small-scales, which play an important role in maintaining the large-scale cir-
 60 culation [Greatbatch *et al.*, 2010a,b; Waterman and Jayne, 2010; Waterman *et al.*, 2011;
 61 Kang and Curchitser, 2015]; with satellite observations only providing surface informa-
 62 tion. 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 sub-surface quantities. Illuminating unresolved quantities using ‘seen’ quantities would extend the reach of existing datasets, 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 [McGovern *et al.*, 2017; Esteves *et al.*, 2018], climate model parameter sensitivity studies [Anderson and Lucas, 2018], chaotic dynamical systems forecasting [Pathak *et al.*, 2018a,b; Vlachas *et al.*, 2018], and parameterising unresolved atmospheric processes [Gentine *et al.*, 2018; Brenowitz and Bretherton, 2018; Jiang *et al.*, 2018; O’Gorman and 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 regards to machine learning, Chapman and Charantonis [2017] constructed a form of neural network known as a self-organising map to reconstruct sub-surface velocities in the Southern ocean using satellite altimetry data and Argo floats. Other studies have used random forests to predict sub-surface 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 sub-surface velocities); as is the case with OGCMs, where high-resolution data is less readily available due to the computational cost. A similar challenge is when data is only available for particular regions, such as mooring data [Hogg, 1992] or gliders [Rudnick *et al.*, 2004; Davis *et al.*, 2008]. 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 a deep neural networks, which are known for a strong ability to generalise [Krizhevsky *et al.*, 2012; LeCun *et al.*, 2015; Goodfellow *et al.*, 2016].

However, deep neural networks are typically considered a ‘black box’, i.e., 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 [Tracey *et al.*, 2015; Kutz, 2017], where the studies of Ling *et al.* [2016a,b] in particular show that a neural network can respect Galilean invariance by utilising the invariant tensors of Pope [1975]. The studies of Ling *et al.* [2016a,b] are important in moving towards 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 sub-surface flow fields, and to determine if they are suited to situations where data is limited to a particular region. To move towards these aims, as a proof-of-concept we will address the following questions:

- 111 1. Can convolutional neural networks represent the spatio-temporal variability of the
112 sub-grid eddy momentum forcing.

- 113 2. How sensitive are the neural networks to the physical processes occurring within
 114 each region, and how well do they generalise to ocean models in different configu-
 115 rations.
 116 3. Is it possible to physically-constrain neural networks to respect global momentum
 117 conservation.
 118 4. Using only information at the surface, can neural networks predict the sub-surface
 119 flow fields.

120 By using data from an idealised high-resolution ocean model, we show that con-
 121 volutional neural networks can represent both the spatial and temporal variability of the
 122 eddy momentum forcing. The region the neural network is trained on, and therefore the
 123 dynamical processes occurring within that region, significantly impact the performance of
 124 the neural network. In particular, training on the most turbulent region produces the best
 125 overall performing neural network. The neural networks successfully generalise to mod-
 126 els with different viscosity coefficients and external wind forcings. Initially momentum is
 127 not conserved globally, but the neural networks can be constrained to respect momentum
 128 conservation without a significant reduction in accuracy. A neural network can accurately
 129 predict the sub-surface flow field when there is a strong barotropic component to the flow.

130 The paper is organised as follows. The quasi-geostrophic ocean model, the degrad-
 131 ing of model data, and convolutional neural network, are introduced in Section 2. Perfor-
 132 mance diagnostics of the neural networks, in terms of non-local predictions and generalis-
 133 ing to different model configurations, are presented in Section 3. We explore methods of
 134 physically-constraining the neural networks in Section 4. Section 5 presents a neural net-
 135 work trained to predict sub-surface flow fields using only information at the surface. We
 136 summarise and discuss our results in Section 6.

137 2 Data and Methods

138 2.1 Quasi-Geostrophic Ocean Model

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

142 The configuration of this model leads to two large-scale circulation gyres separated
 143 latitudinally by a strong meandering zonal jet. The model is configured to represent an
 144 idealised version of current systems such as the Gulf Stream in the North Atlantic or the
 145 Kuroshio Extension in the North Pacific; both these current systems exhibit vigorous ed-
 146 dies interacting with a strong mean-flow. The time-mean streamfunction, which illustrates
 147 the double-gyre flow structure, can be seen in Figure 1a of Mana and Zanna [2014].

148 The potential vorticity q is given by

$$149 \quad q = \nabla^2 q + \beta y + \frac{\partial}{\partial z} \left(\frac{f_0^2}{N^2} \frac{\partial \psi}{\partial z} \right), \quad (1)$$

150 where $f = f_0 + \beta y$ is the planetary vorticity, f_0 is the Coriolis parameter, $\beta = df/dy$
 151 is the Rossby parameter, $\nabla = (\partial/\partial x, \partial/\partial y)$ is the horizontal gradient operator, $N =$
 152 $(-\frac{g}{\rho} \frac{d\rho}{dz})^{1/2}$ is the Brunt-Väisälä frequency, g is gravity, ρ is density, and ψ is the stream-
 153 function for the non-divergent horizontal velocity $\mathbf{u} = (-\partial\psi/\partial y, \partial\psi/\partial x)$.

154 The model has three layers ($m = 1$ upper, $m = 2$ middle, $m = 3$ upper), with thick-
 155 nesses H_m of 250 m, 750 m, 3000 m, respectively. For each layer, the following prognos-
 156 tic equation is solved

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

where $\mathcal{D} = \nu \nabla^4 \psi - r \nabla^2 \psi \delta_{m,3}$ 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 ν . The second dissipation term parameterises the presence of an Ekman 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

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

where $g(x) = L/2 + 0.2(x - L/2)$, $L = 3840$ km is the domain length, and ρ_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 *Manai and Zanna* [2014]; *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 2D low-pass filter, in order to produce data that is 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 sub-filter 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 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 (sub-filter) variables onto the resolved (filtered) variables. Ignoring vertical effects and planetary vorticity, the horizontal momentum equation is given by

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

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

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

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

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

190 The low-pass filtering operation results in an additional forcing term in Equation 7
 191 for the filtered momentum; the additional momentum forcing \mathbf{S} is given by Equation 8,
 192 the divergence of a Reynolds stress. The vector $\mathbf{S} = (S_x, S_y)$ represents the effects of the
 193 sub-filter momentum field on the filtered momentum field, i.e., the interaction between
 194 small-scale eddies and the large-scale flow. As the sub-filter eddy forcing \mathbf{S} depends on
 195 the sub-filter variables, it requires a physical parameterisation or closure.

196 2.3 Predictive Algorithm: Convolutional Neural Networks

197 Convolutional Neural Networks (CNNs) have proven successful in many areas of
 198 computer vision [Krizhevsky *et al.*, 2012; Simonyan and Zisserman, 2014; Dong *et al.*,
 199 2016], where the primary objective is to extract information from an image, in order to
 200 perform a particular task. CNNs work by applying successive layers of convolutions (a
 201 form of spatial filtering) to the input; the complexity of the extracted information increases
 202 with the number of convolution layers. The powerful property of CNNs is that the filters
 203 of each convolution are learnt as part of the training process - they are not specified a pri-
 204 ori. Therefore CNNs learn to extract the most ‘useful’ information from the input variable,
 205 given training on a particular dataset.

206 We chose to use CNNs, as opposed to a deep neural network of multiple fully-
 207 connected layers, due to their superior performance in computer vision tasks where the
 208 inputs have a two dimensional structure [Krizhevsky *et al.*, 2012]. We wanted a machine
 209 learning algorithm that could exploit the two dimensional lateral structure of turbulent
 210 fluids. Spatial filtering of the equations of motion of turbulent fluids is not new, and is
 211 used in Large Eddy Simulation (LES) [Moeng, 1984; Sagaut, 2006]. Therefore, the learnt-
 212 filtering operations of a CNN appeared to be a natural choice of data-driven algorithm to
 213 apply to geophysical flows.

214 The training process involves the minimisation of an appropriately defined loss func-
 215 tion, which measures the difference between the output of the CNN, and the desired tar-
 216 gets. If the optimisation procedure was successful, such that the loss function on previ-
 217 ously unseen data converges, the CNN will have learnt to extract the most important infor-
 218 mation from the input. The CNN then uses the information to predict continuous values.
 219 The CNN constructs the final prediction through a linear regression layer, which regresses
 220 the desired output onto the final feature maps (feature maps are the intermediate results of
 221 each convolution layer).

222 Here we use CNNs to represent the sub-filter eddy momentum forcing. The input
 223 is the filtered-streamfunction ψ of the upper vertical layer, which represents our resolved
 224 variable that the neural networks will extract information from. The output variables are
 225 the zonal S_x and meridional S_y components of the sub-filter momentum forcing \mathbf{S} , defined
 226 by Equation 8. An example input and output is shown in Figure 1. Separate CNNs are
 227 trained for each component of the sub-filter momentum forcing S_x and S_y . We only con-

228 sider data from the upper-layer of the model; this is because the flow is surface-intensified,
 229 and we are assuming that our filtered quantities are similar to satellite altimetry data,
 230 which only provide information at the surface.

231 In addition to testing whether it is possible to train a neural network to predict S_x
 232 and S_y , from $\bar{\psi}$, we explore how a neural network trained on one region performs on an-
 233 other previously unseen region, i.e. how important local vs non-local information is for
 234 different regions. We therefore construct three different datasets from the QG model data,
 235 one for each region being studied. We choose regions which differ most in their dynamical
 236 behaviour, and are shown in Figure 1a: Region 1 is near the jet-separation point of the
 237 western boundary, where there is a strong, inertial zonal jet. Region 2 is near the eastern
 238 boundary downstream of the jet extension, where the dynamics are more wave-like in na-
 239 ture. Region 3 is in the centre of the southern gyre, which is energetically less active than
 240 regions 1 and 2.

241 Data from the three regions are split temporally into training and validation datasets.
 242 The 10-years of daily data (3650 days) are split into the first ~9 years (3300 days) to train
 243 the neural networks, and the final year (350 days) is set aside for validation. To reduce
 244 the computational cost, and the number of parameters of each CNN, we split each region
 245 spatially from the initial 160×160 grid points, to sixteen 40×40 grid point sub-regions,
 246 as depicted in Figure 1c. Reducing the input and output size of the neural network from
 247 160×160 to 40×40 significantly decreases the number of trainable weights, and therefore
 248 the computational cost (we attempted to make predictions for the full 160×160 of each
 249 training region, but this led to a neural network with over 250,000,000 parameters, which
 250 was computationally impractical).

251 Making predictions for a 40×40 area instead of a 160×160 area also increases the
 252 amount of training and validation data by a factor of sixteen, from 3300 and 350 samples,
 253 to 52800 and 5600 respectively, where a sample is defined as a single input-output pair
 254 of the neural network. We therefore have 52800 spatial maps (size 40×40 grid points) of
 255 input-output pairs to train the neural networks, and 5600 spatial maps of input-output pairs
 256 set aside for validation.

257 We train CNNs to separately predict S_x and S_y , using data from three different re-
 258 gions of the model; this gives a total of 6 neural networks. Each neural network is de-
 259 noted 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
 260 the trained weights of the neural network, and $R = 1, 2, 3$ refers to the region on which the
 261 neural network has been trained. For example, the neural network trained on region 2 to
 262 predict the meridional component S_y is denoted by $f_y(\bar{\psi}, \mathbf{w}_2)$.

263 To distinguish predictions from the true values, we label neural network predictions
 264 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 sub-filter momentum
 265 forcing remain as S_x, S_y . We use the mean-squared error as the loss function,

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

266 which quantifies the difference between the neural network predictions and the truth, and
 267 where the summation is over all samples. The neural networks are trained (i.e. optimised)
 268 using a form of stochastic gradient descent, namely the Adam optimisation algorithm
 269 [Kingma and Ba, 2014], which minimises the loss function L defined in Equation 9. The
 270 training of each neural network $f_i(\bar{\psi}, \mathbf{w}_R)$, iteratively adjusts the values of the weights \mathbf{w}_R ,
 271 such that the loss function in Equation 9 is minimised. Therefore each neural network has
 272 a different set of weights \mathbf{w}_R ; it is these weights which determine how each neural net-
 273 work extracts information and makes predictions.

274 The architecture used for each $f_i(\bar{\psi}, \mathbf{w}_R)$ contains three convolution layers, a max
 275 pooling layer, and a final fully-connected layer (Figure 1). The max pooling layer reduces
 276 the dimensionality of the previous layer, by selecting the maximum value within a 2×2

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 non-linear 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 towards zero mean and unit variance, removing the need for batch normalisation - batch normalisation 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 minimises 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 *et al.*, 2015], with the Tensorflow backend [Abadi *et al.*, 2016]. Before training, all datasets are separately normalised 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 optimisation 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 $\bar{\psi}$ from the validation dataset, i.e., the final year of withheld data. We make predictions for the full-domain to determine how each neural network generalises to unseen, dynamically-distinct, regions. As the input and output size of each neural network is 40×40 grid points, we tile together predictions for the full domain of size 512×512 ; 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 sub-surface flow field, using only information at the surface, we train a new neural networks. The new neural network has an identical architecture to those discussed previously, and is trained to predict the middle-layer streamfunction using the upper-layer streamfunction as the input; this neural network is described in more detail in Section 5.

3 Neural Network Generalisation and Sensitivity

3.1 Non-Local Predictions

The filtered streamfunction represents for example observational measurements from satellite altimetry or coarse-resolution model data. The sub-filter eddy momentum forcing represents unresolved turbulent processes. Our goal is to replicate the complex spatio-temporal variability of S_x and S_y using neural networks $f_i(\bar{\psi}, \mathbf{w}_R)$. However observational data such as moorings [Hogg, 1992] or gliders [Rudnick *et al.*, 2004; Davis *et al.*, 2008], 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 generalise to previously un-seen regions.

We study the spatio-temporal 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 512×512 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 (Figure 2a, e, and i). In particular, strong meanders which extend eastward from the western boundary are visible. The amplitude of the spatio-temporal variability of S_x ($1.4 \times 10^{-6} \text{ms}^{-2}$) is of similar magnitude to the time-mean ($1.5 \times 10^{-6} \text{ms}^{-2}$).

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 (Figure 2b, c, and d). Their magnitudes however vary significantly. The predictions of $f_x(\bar{\psi}, \mathbf{w}_1)$, trained on data from the western boundary, are almost identical to the true S_x , and successfully reproduces the correct amplitude and variability (Figure 2b, f, j). The neural network $f_x(\bar{\psi}, \mathbf{w}_2)$, 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 $f_x(\bar{\psi}, \mathbf{w}_3)$, trained on the southern gyre, underestimates the true S_x by an order of magnitude (Figure 2d, h, l).

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 $f_x(\bar{\psi}, \mathbf{w}_1)$ and $f_x(\bar{\psi}, \mathbf{w}_2)$ are highly correlated with the truth ($r > 0.9$) within the jet, but tend towards zero or negative correlation near the eastern boundary (Figure 2m and 2n). The predictions of $f_x(\bar{\psi}, \mathbf{w}_3)$ have a more consistent positive correlation across the gyres and other more quiescent regions, (Figure 2o).

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 ($1 \times 10^{-6} \text{ ms}^{-2}$) 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, $f_y(\bar{\psi}, \mathbf{w}_1)$, most effectively reproduces the true S_y . However, the time-mean of $f_y(\bar{\psi}, \mathbf{w}_1)$ (Figure 3f) has a positive bias everywhere in the domain, whereas the time-means of $f_y(\bar{\psi}, \mathbf{w}_2)$ and $f_y(\bar{\psi}, \mathbf{w}_3)$ (Figure 3g and 3h respectively) do not.

The correlations between S_y and \tilde{S}_y are similar to the zonal component: $f_y(\bar{\psi}, \mathbf{w}_1)$ and $f_y(\bar{\psi}, \mathbf{w}_2)$ are highly correlated ($r > 0.8$) within the jet, but not in the gyres. In contrast, $f_y(\bar{\psi}, \mathbf{w}_3)$ has a consistently positive correlation across the full domain, despite failing to reproduce the amplitude within the jet. In fact, the correlation of $f_y(\bar{\psi}, \mathbf{w}_3)$ 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 sub-filter 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 behaviour 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 generalising 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

375 regions can lead to significant net biases in the predictions, as seen in Figure 3f. How to
 376 correct for such biases will be discussed in Section 4.

377 3.2 Generalising to Different Reynolds Numbers

378 In Section 3.1, we investigated how neural networks trained on different regions
 379 of the domain generalise to other previously unseen regions. We now test how the neu-
 380 ral networks generalise to different regimes, in particular different Reynolds number. In
 381 Section 3.1, we found that the neural networks trained on region 1, the western bound-
 382 ary, successfully generalised to different regions; we therefore apply $f_x(\bar{\psi}, \mathbf{w}_1)$ to new
 383 model data with different wind stress amplitudes and viscosity coefficients to test its per-
 384 formance. We use models with higher and lower wind forcings, to test regimes which are
 385 both more and less turbulent than the original model, which had a wind stress amplitude
 386 of $\tau_0 = 0.8 \text{ Nm}^{-2}$ and viscosity $\nu = 75 \text{ m}^2\text{s}^{-2}$.

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

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

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

404 3.3 Sensitivity of Neural Networks to Under-Sampling

405 We have so far trained the neural networks with densely sampled data, i.e., we have
 406 data at each grid point for both the input and output variables. However, most observa-
 407 tional datasets are spatially sparse, e.g. Argo floats [Roemmich *et al.*, 2009]. We therefore
 408 explore the impact of under-sampling with a new collection of neural networks trained
 409 on region 1 to predict S_x , but with the training data sub-sampled. At each time-slice of
 410 the training data, we randomly sample (without replacement) N points of the 40×40 in-
 411 put variables, $\bar{\psi}$, and output variables S_x . Using these N randomly sampled values, we
 412 use a cubic interpolation to reconstruct the full 40×40 grid point input and output (with
 413 a nearest-neighbour interpolation for grid points that fall outside the convex hull of the
 414 cubic interpolation).

415 These reconstructed time-slices from sub-sampled data are used to train a new set
 416 of neural networks. We vary the number of points N sub-sampled from $> 90\%$ to $< 5\%$
 417 of the original 1600 points of the input and output variables. We have a neural network
 418 for each value of N , the sub-sampling rate. Using the neural networks trained on under-
 419 sampled data, we calculate the root-mean square error (RMSE) on the final year of val-
 420 idation data over the entire domain. The validation data is not sub-sampled, providing a
 421 stronger and more accurate test of the neural network's performance.

The RMSE is shown as a function of percentage of points sampled (Figure 7c). We find that the RMSE increases significantly only when the percentage of spatial points sampled drops below 10% (the error doubles at a sub-sampling 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 non-linear interpolation. The spatial map of RMSE of the neural network trained with 18.75% sub-sampled data (Figure 7b) shows minimal changes relative to the neural network trained on the original (un-altered) training data (Figure 7a). 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 under-sampling, where the 40×40 input and output grid of the neural network is spaced out over the entire domain. In other words, we sub-sample the input and output variables of the original 512×512 grid to a regularly spaced 40×40 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 streamfunction (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 sub-sampling the input streamfunction, the local information relevant to estimate velocities and velocity shears is lost.

Therefore the success of training convolutional neural networks on observational datasets 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 7. We anticipate that the use of (interpolated) datasets with horizontal resolution of 1 degree 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 datasets 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 $f_i(\psi, \mathbf{w}_R)$, and then explore different methods of imposing conservation of momentum globally.

4.1 Momentum Biases

Each sub-region (including those used to train the neural networks) may have a non-zero spatially-integrated momentum tendency. However, globally, the true sub-filter momentum forcing \mathbf{S} should re-distribute momentum, and not act as a source or sink, i.e. $\iint S dx dy = 0$. 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 sub-region, we expect to have imperfect momentum conservation, which will depend upon the particular dynamical processes within each region. For example, if eddies within a particular 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 generalise the (local) input of momentum to the rest of the domain. A net source or sink of momentum will manifest as a non-zero 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 (Figure 5a and 5b), 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 (Figure 5c and 5d respectively) over the full domain, we observe significant non-zero biases.

Consider the zonal component of the eddy momentum forcing in Figure 5c: $f_x(\bar{\psi}, \mathbf{w}_1)$ 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 $0.0008 (10^{-6} \text{ms}^{-2})$ respectively; this leads to zonal velocity changes of $\Delta u = 0.95, 0.63$, and $0.025 (\text{ms}^{-1})$. These changes are of similar magnitude to the time-mean zonal flow, which peaks at approximately 0.9 ms^{-1} within the jet core.

There are also significant biases in the predictions of the meridional component \tilde{S}_y , shown in Figure 5d. The positive bias of $f_y(\bar{\psi}, \mathbf{w}_1)$ 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 $\langle \tilde{S}_y \rangle = 0.02, -0.01$, and $0.002 (10^{-6} \text{ms}^{-2})$ for $f_y(\bar{\psi}, \mathbf{w}_1)$, $f_y(\bar{\psi}, \mathbf{w}_2)$, and $f_y(\bar{\psi}, \mathbf{w}_3)$ respectively, leads to the following changes: $\Delta v = 0.63, -0.31$, and $0.06 (\text{ms}^{-1})$. Some of these changes are the same magnitude as the time-mean meridional flow.

4.2 Towards Momentum-Conserving Neural Networks

The predictions of neural networks $f_x(\bar{\psi}, \mathbf{w}_1)$ and $f_y(\bar{\psi}, \mathbf{w}_1)$, 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 non-zero 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 trial three approaches (A, B, and C) to reduce the biases identified in Figure 5c and 5d.

- (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 sub-filter momentum forcing, but with momentum conservation intrinsically embedded. I.e. same training data, but altered architecture. The motivation behind this approach is that if the local source of momentum within the 40×40 output grid is zero, then this may reduce the global net source of momentum.
- (B) Pre-processing of input: Train on region 1 with the original architecture described in Table 1 but with the spatial-mean removed from each 40×40 training snapshot of S_x . In other words, remove the local bias from each training snapshot. If the local source of momentum of each 40×40 sample is zero within the training data, then training on such data may reduce local biases when making future predictions on un-seen data. However, does not guarantee that subsequent predictions will have zero local bias. This approach tries to increase the probability of *learning* local momentum conservation.
- (C) Post-processing of output: train on region 1, and enforce global momentum conservation after the predictions have been made. I.e. no changes to training data or architecture, but with additional processing of the full-domain predictions \tilde{S}_x and \tilde{S}_y .

518 The associated neural networks of each approach are labelled as $f_i(\bar{\psi}, \mathbf{w}_1^A)$, $f_i(\bar{\psi}, \mathbf{w}_1^B)$, and
 519 $f_i(\bar{\psi}, \mathbf{w}_1^C)$ respectively, where $i = (x, y)$ denotes either the zonal S_x or meridional S_y com-
 520 ponent being predicted.

521 All neural networks are optimised using the same training parameters given in Ta-
 522 ble 1. Approach A, which alters the architecture, and approach B, which alters the train-
 523 ing data, are enforcing momentum conservation not just globally, but within the 40×40
 524 sub-region being predicted. This local conservation is useful for enforcing global conser-
 525 vation. However local conservation may not be desirable if there's convergence of eddy
 526 momentum fluxes in a particular region, which can impact the large-scale flow, e.g. if ed-
 527 dies are fluxing momentum into the jet at a particular along-stream position, enforcing lo-
 528 cal conservation in a neural network may lead to missing these effects. Therefore caution
 529 must be taken with restricting architectures in this way.

530 We now explore the performance of the newly constrained neural networks and the
 531 net momentum input relative to that of the original neural networks trained on region 1:
 532 $f_x(\bar{\psi}, \mathbf{w}_1)$ and $f_y(\bar{\psi}, \mathbf{w}_1)$. The spatial-averages of neural networks based on approaches A,
 533 B, and C are shown in Figure 6, with the same scale axes as in Figure 5.

534 Approach B has significant biases of approximately -0.01 and -0.015 (10^{-6}ms^{-2}) in
 535 the zonal and meridional components respectively; the optimisation procedure aims to re-
 536 produce the *variability* in the training data, and not spatial-means, therefore pre-processing
 537 the training data does not remove the biases. Compared to the original neural networks
 538 trained on region 1, the biases of approaches A and C are 3 to 5 orders of magnitude
 539 lower, in both the zonal and meridional components. The post-processing approach is ex-
 540 actly zero by construction, while the altered-architecture approach A is not exactly zero
 541 due to the overlapping-tiling procedure. The biases of $f_x(\bar{\psi}, \mathbf{w}_1^A)$ and $f_y(\bar{\psi}, \mathbf{w}_1^A)$ are ap-
 542 proximately -0.002 and -0.0005 (10^{-6}ms^{-2}) which, over the course of a year, would lead
 543 to velocity changes of $\Delta u = -0.06$ and $\Delta v = -0.01$ (ms^{-1}) respectively - now an order of
 544 magnitude smaller than the time-mean flow.

545 The correlation maps of all momentum-conserving approaches (not shown) change
 546 little from the original correlation maps of $f_x(\bar{\psi}, \mathbf{w}_1)$ and $f_y(\bar{\psi}, \mathbf{w}_1)$, shown in Figure 2m and 3m
 547 respectively. All approaches reproduce the correct spatial patterns of the true S_x and S_y
 548 (e.g., Figure 6 for standard deviations). However, approaches A and B underestimate the
 549 amplitude of S_x and S_y by approximately 20-30%, whereas there is a little difference be-
 550 tween approach C and the truth (< 10%).

551 In summary, approach C of post-processing successfully enforces momentum conser-
 552 vation, without sacrificing accuracy in the predictions of the eddy momentum momentum
 553 forcing. Approach B, altering the training data, was not efficacious at reducing the net bi-
 554 ases. The physically-constrained architecture of approach A successfully reduced the net
 555 bias, but at the expense of 20-30% accuracy. Though further altering of the architecture
 556 (e.g. increasing number of convolution layers and filters) or training procedure (decreas-
 557 ing the learning rate, with increased number of training epochs) could reduce this drop in
 558 accuracy by countering the restriction placed on the architecture.

5 Predicting Sub-Surface Flow

560 We have shown that neural networks, by using the filtered-streamfunction as the in-
 561 put variable, can provide information on unresolved turbulent processes, namely the sub-
 562 filter momentum forcing. We have assumed that the filtered-streamfunction represents
 563 some limited set of observations, or data from a coarse-resolution ocean model. However,
 564 coarse-resolution ocean models still produce data for below the surface, whereas satellite
 565 observations do not. Here we address the issue of inferring sub-surface information solely
 566 from surface fields. Our approach is conceptually similar to *Chapman and Charantonis*
 567 [2017], which used a form of neural network called a self-organising map to reconstruct

568 sub-surface velocities in the Southern ocean, using satellite altimetry and Argo float data.
 569 Using the QG model data described in Section 2.1, we test whether a neural network can
 570 predict the middle-layer streamfunction, using only the surface filtered-streamfunction.

571 We train a new neural network $\tilde{\psi}_2 = f(\bar{\psi}_1, \mathbf{W})$ (which has the same architecture as
 572 before, but with a different output and weights) to minimise the mean-squared error loss
 573 function $L \propto (\psi_2 - \tilde{\psi}_2)^2$, where $\bar{\psi}_1$ is the filtered-streamfunction of the upper-layer, ψ_2
 574 is the true streamfunction of the middle-layer, and $\tilde{\psi}_2$ is the neural network predictions.
 575 Again, to assess the ability to generalise to unseen regions, we only train the neural net-
 576 work on the western boundary (training region 1). Diagnostics of the true ψ_2 and predic-
 577 tions $\tilde{\psi}_2$, including the correlation between them, are shown in Figure 8a-e. The neural
 578 network accurately reproduces the middle-layer time-mean and standard deviation of the
 579 streamfunction within the jet region. The neural network accurately reproduces the correct
 580 amplitude of the true ψ_2 within the jet, but underestimates the amplitude by $\approx 50\%$ within
 581 the gyres. Independent of the amplitude, the predictions $\tilde{\psi}_2$ are highly correlated ($r > 0.8$)
 582 almost everywhere in the domain with the true ψ_2 .

583 The decrease in accuracy in the gyres is likely due to only training within the west-
 584 ern boundary, where the streamfunctions of the upper- and middle-layers are more tightly
 585 coupled due to the strong barotropic nature of the flow. Within the gyres, the barotropic
 586 component is not as dominant - this could cause the neural networks to underestimate the
 587 amplitude away from the jet. Alternatively the adjustment time scales of the upper- and
 588 middle-layers are not the same, which perhaps requires more training data in order to cap-
 589 ture interactions over longer time scales.

590 We take the approach one step further, by predicting the bottom-layer streamfunc-
 591 tion, using the same neural network and its weights $f(\bar{\psi}_1, \mathbf{W})$, but now using the predic-
 592 tions of the middle-layer streamfunction as the input, i.e., $\tilde{\psi}_3 = f(\bar{\psi}_2, \mathbf{W})$. We test whether
 593 a neural network trained to predict the middle-layer streamfunction can provide any infor-
 594 mation on the bottom-layer streamfunction (without re-training), by inputting the middle-
 595 layer streamfunction as an input. Mathematically, this is written as $\tilde{\psi}_3 = f(f(\bar{\psi}_1, \mathbf{W}), \mathbf{W})$.

596 Diagnostics of the true (ψ_3) and predicted ($\tilde{\psi}_3$) bottom-layer streamfunction are
 597 shown in Figure 8f-j. Despite a moderate correlation of $r \approx 0.5$ across the domain, the
 598 predictions fail to reproduce the correct time-mean, which has a circulation in the oppo-
 599 site direction to the truth. This is due to the neural network being trained to predict the
 600 middle-layer flow, which on average is more aligned with the upper-layer. Therefore when
 601 the neural network is given the middle-layer streamfunction as an input, it predicts the
 602 bottom-layer flow as on-average being in the same direction, which is not the case. The
 603 neural network also hasn't be trained to predict the effects of the additional bottom drag,
 604 decreasing the accuracy further - more data could improve this issue, as the longer time
 605 scales associated with bottom drag may be absent from the training dataset.

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

612 6 Conclusions & Discussion

613 6.1 Summary

614 In this study, we have demonstrated as a proof-of-concept that machine learning al-
 615 gorithms can provide information on unresolved turbulent processes, when given a smoothed-
 616 view of the dynamics (i.e. the filtered-streamfunction). 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, i.e. 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 is limited to only particular regions, for example mooring data [Hogg, 1992] or gliders [Rudnick *et al.*, 2004; Davis *et al.*, 2008]. Training on a sub-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 generalisation, 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 idealised 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 down-stream the eddies drive the time-mean circulation. Therefore training neural networks on different along-stream positions may lead to different dynamical-processes being learnt, 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- e.g. 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 generalises exceedingly well to QG models with different viscosity coefficients and wind forcings (Figure 4). The neural network within the jet reproduces the correct spatio-temporal 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 5c and 5d), we show that momentum conservation can be enforced without a significant reduction in accuracy (Figure 6), through either a physically-constrained architecture or post-processing of the predictions.

We also show that a new neural network can be trained to predict the middle-layer streamfunction, using only the upper-layer streamfunction as the input, i.e., 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 streamfunction to predict the bottom-layer streamfunction, the neural network captures some of the variability, but fails to replicate the time-mean of the true bottom-layer streamfunction ψ_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 datasets [Chapman and Cha-

668 *rantonis*, 2017; *Su et al.*, 2018; *Giglio et al.*, 2018], the present work demonstrates that
 669 convolutional neural networks in particular are an excellent tool for such tasks. Neural
 670 networks should be further tested and exploited in the future for data inference due to

- 671 • their resilience, such that accurate predictions for the full domain can be generated
 672 by training on a sub-region;
 673 • their generalisation to different external forcings, without any further training such
 674 that predictions outside the regime trained on can be successful;
 675 • their ability to be successfully trained with under-sampled data (Figure 7).

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

686 6.3 Implications for parametrizations

687 Although we have motivated our study through the leverage of observations and
 688 coarse-resolution model data, our results have implications for eddy parameterisations of
 689 momentum, and more generally for sub-grid parametrizations. As discussed previously,
 690 machine learning has been used to parameterise unresolved processes in the atmosphere
 691 [*Brenowitz and Bretherton*, 2018; *Jiang et al.*, 2018; *Gentine et al.*, 2018; *O’Gorman and*
 692 *Dwyer*, 2018].

693 We have shown that neural networks can successfully represent the spatio-temporal
 694 variability of the eddy momentum forcing, implying potential for data-driven oceanic tur-
 695 bulence closures in the future, as suggested by *Zanna et al.* [2018]. The generalisation
 696 ability of the neural networks shows that only a limited amount of observations or high-
 697 resolution model data may be needed to successfully represent sub-grid scale processes.
 698 While the CNNs are successful at representing relationship between the eddy momentum
 699 forcing and their effect on the resolved flow, the low-resolution climate models might have
 700 biases that are too severe (e.g., weak transport and velocity shears) to lead to a successful
 701 representation of the eddy momentum forcing from CNNs as trained here. Nonetheless,
 702 our results also show that they perform very well for different amplitude of forcing and
 703 dissipation. Therefore, until the CNNs are implemented into a coarse-resolution ocean
 704 model, their success in improving numerical simulations is speculative but deserves to be
 705 investigated.

706 Whether neural networks are being used to leverage observations, or more impor-
 707 tantly to construct a data-driven eddy parameterisation, caution must be taken to ensure
 708 that the laws of physics are respected. More work into physically-constrained machine
 709 learning algorithms is crucial, and successful applications of data-driven techniques should
 710 incorporate physical knowledge. Indeed, the neural network turbulence model of *Ling*
 711 *et al.* [2016b] out-performed more simple linear models only when Galilean invariant
 712 stress tensors from *Pope* [1975] were used, which are also a key ingredient of the eddy
 713 parameterisation proposed by *Anstey and Zanna* [2017]. As previously discussed, we suc-
 714 cessfully enforce global momentum conservation in the present work, such that future im-
 715 plementations of data-driven parameterisations, despite being semi-empirical, can be al-
 716 tered to respect physical principles. Specifically, physical constraints can be incorporated
 717 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 and Sun, 2015] is given by $O(\sum_l^d n_{l-1} \cdot s_l^2 \cdot n_l \cdot m_l^2)$, 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 [Howard *et al.*, 2017, e.g.], 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, i.e. make predictions for a region smaller than 40×40 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 parametrizations, the training of the neural networks is only done once.

6.4 Future Work

Our study is a step towards using convolutional neural networks to extend the reach of currently available observational or model data. Our proof-of-concept study was conducted in an idealised QG model. The next stage involves training neural networks on actual observational datasets (as described in Section 7.2) or on more realistic model data (e.g. a 1/40th degree global model which resolves the mesoscale and submesoscale eddy fields, such as in Rocha *et al.* [2016]).

We used nine years of data to train the neural networks, and one year for validation. Gentine *et al.* [2018] showed, with regards to parameterising convection with neural networks, that the training dataset could be reduced in size from 12 months to 3 months, with little change in the overall mean-squared 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 under-sampling 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).

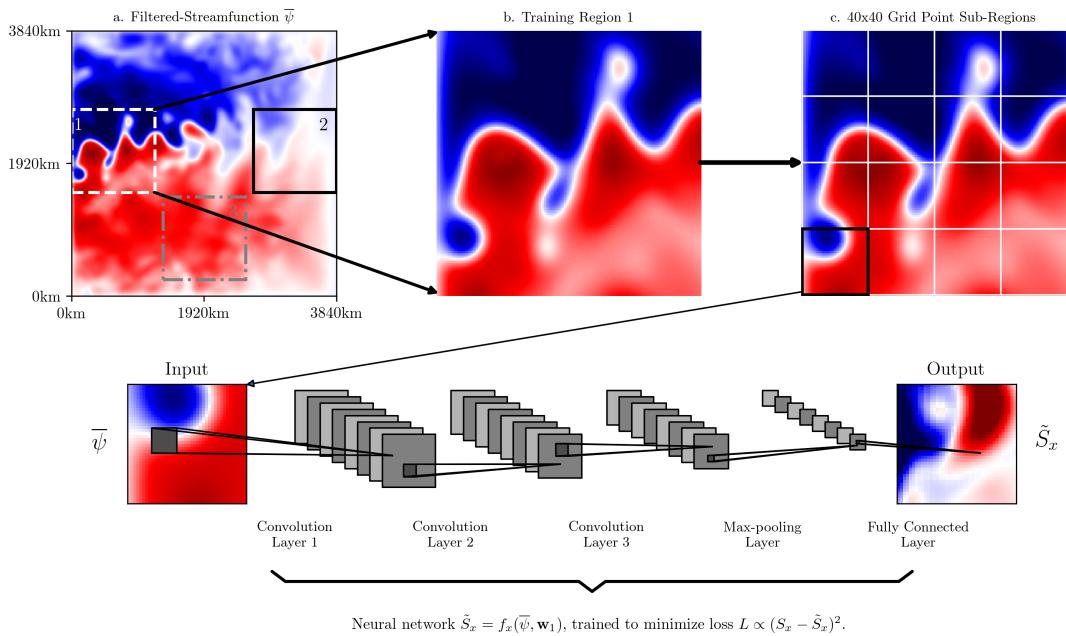
Training on the western boundary produces the best performance. However, the high skill within the jet does not fully translate to high skill in all parts of the gyres. The best correlations in the gyres occurs instead when training on the southern gyre, and not the western or eastern boundaries (Figure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

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

where the summation is over all regions, and \tilde{S}_x^{OPT} is the corresponding optimal prediction (with an analogous \tilde{S}_y^{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 specialisation, and the skill reduces as the single neural network simply ‘averages’ the effects of the three regions

765 together. We will attempt to implement the neural networks (as trained here, or as a com-
766 bination of neural networks) into a coarse resolution version of the QG model to test their
767 performance as a sub-grid scale parametrization.

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



772
 773 **Figure 1.** Panel (a) illustrates the upper-layer filtered-streamfunction $\bar{\psi}$ of the QG model, including the
 774 three regions in which we train the neural networks: region 1 (white-dashed) is on the western boundary,
 775 region 2 (black-solid) is on the eastern boundary, and region 3 (grey-dash-dotted) is centered on the southern
 776 gyre. Panel (b) shows a close-up of the filtered-streamfunction $\bar{\psi}$ within training region 1 while Panel (c)
 777 illustrates how training region 1 is split into 16 40×40 grid point sub-regions - the size of the input and output
 778 arrays of the neural network is 40×40 grid points. The input variable of each neural network is the filtered
 779 streamfunction $\bar{\psi}$, and the output variable is either the zonal component \tilde{S}_x or meridional component \tilde{S}_y of
 780 the sub-filter eddy momentum forcing. The architecture of the convolutional neural network, with an example
 781 input $\bar{\psi}$ and output \tilde{S}_x , is illustrated underneath Panels (a), (b), and (c).

819

Acknowledgments

820
 821
 822
 823
 824
 825
 826

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 which 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>.

827

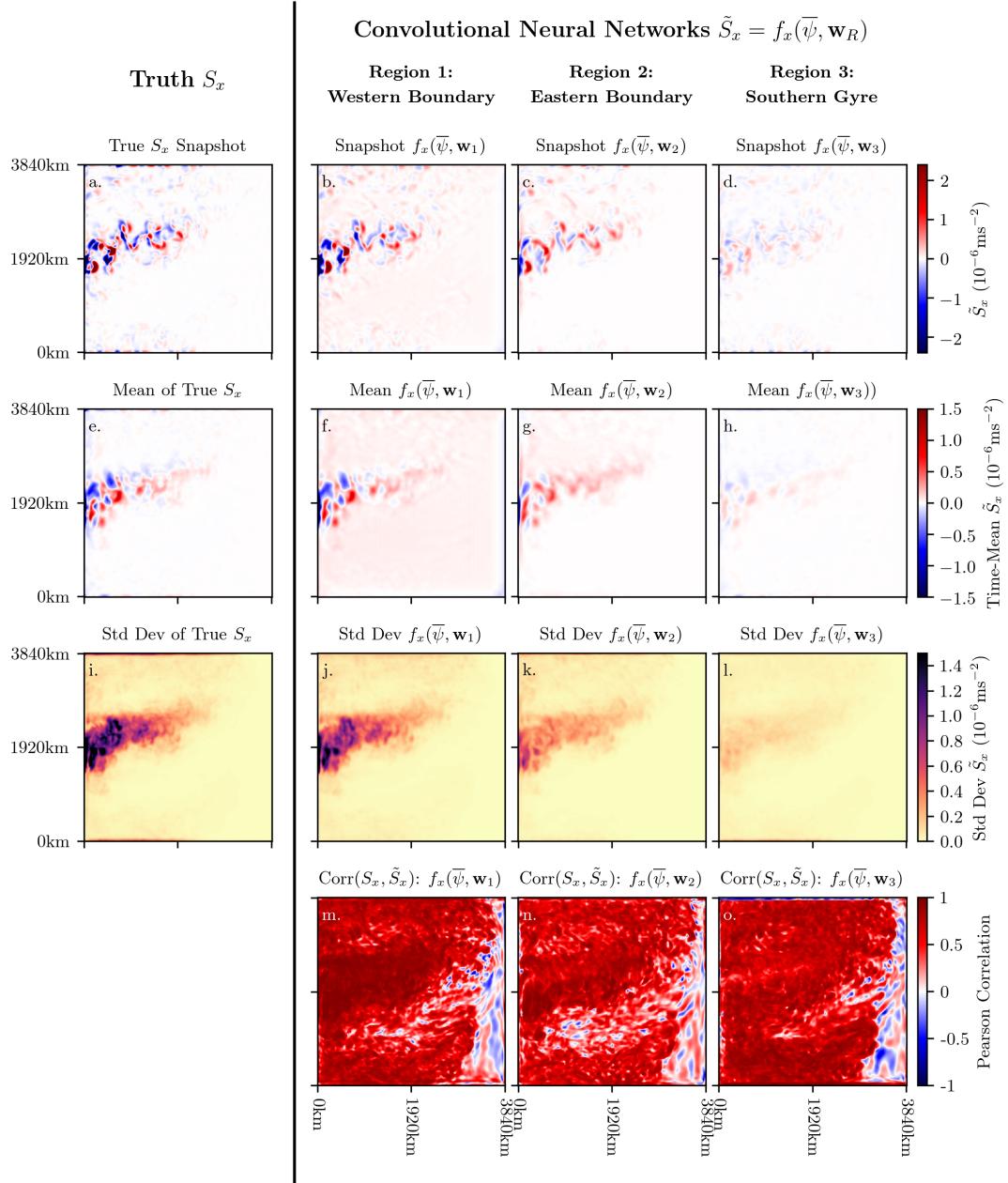
References

828
 829
 830
 831
 832
 833
 834
 835
 836

- Abadi, M., P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al. (2016), Tensorflow: A system for large-scale machine learning., in *OSDI*, vol. 16, pp. 265–283.
- Abernathey, R. P., and J. Marshall (2013), Global surface eddy diffusivities derived from satellite altimetry, *Journal of Geophysical Research: Oceans*, 118(2), 901–916.
- Anderson, G. J., and D. D. Lucas (2018), Machine learning predictions of a multiresolution climate model ensemble, *Geophysical Research Letters*, 45(9), 4273–4280.
- Anstey, J. A., and L. Zanna (2017), A deformation-based parametrization of ocean mesoscale eddy reynolds stresses, *Ocean Modelling*, 112, 99–111.

Quasi-Geostrophic Model Parameters	
Domain size (grid points)	512×512
Domain length (L)	3840 km
Resolution (Δx)	7.5 km
Viscosity (ν)	75 m ² s ⁻¹
Rossby deformation radii (L_{Ro})	40.23 km
Velocity scale (\sqrt{EKE})	0.21 ms ⁻¹
Planetary vorticity (f_0)	10 ⁻⁴ s ⁻¹
Rossby parameter (β)	2 * 10 ⁻¹¹ m ⁻¹ s ⁻¹
Gravity (g)	9.8 ms ⁻²
Reduced gravity (g')	0.034, 0.018 ms ⁻²
Bottom drag coefficient (r)	4 * 10 ⁻⁸ s ⁻¹
Wind stress amplitude (τ_0)	0.8 Nm ⁻²
Reference density (ρ_0)	10 ³ kgm ⁻³
Neural Network Data Details	
Data source	Quasi-geostrophic ocean model
Input variable (feature)	Filtered streamfunction $\bar{\psi}$
Output variables (targets)	Sub-filter 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	52800 (years 1-9)
Number of validation samples	5600 (year 10)
Standardisation 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
Optimiser	Adam
Learning rate	0.001
Momentum	0.9
Batch size	16
Training epochs	200

781 **Table 1.** Details on the following: the quasi-geostrophic ocean model parameters, the datasets used to train
 782 the neural networks, the architecture parameters, and the optimisation parameters.



783 **Figure 2.** Examining the non-local prediction ability. Comparisons of the true zonal component of the
 784 sub-filter momentum forcing S_x , with the neural networks trained using data from three different regions. The
 785 first three rows compare snapshots, time-means, and the standard deviation respectively, while the bottom row
 786 shows the correlation between the true S_x and the predictions \tilde{S}_x . The first column contains the diagnostics
 787 using the true zonal sub-filter momentum forcing S_x , while columns two, three, and four use predictions \tilde{S}_x
 788 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
 789 using the validation data.

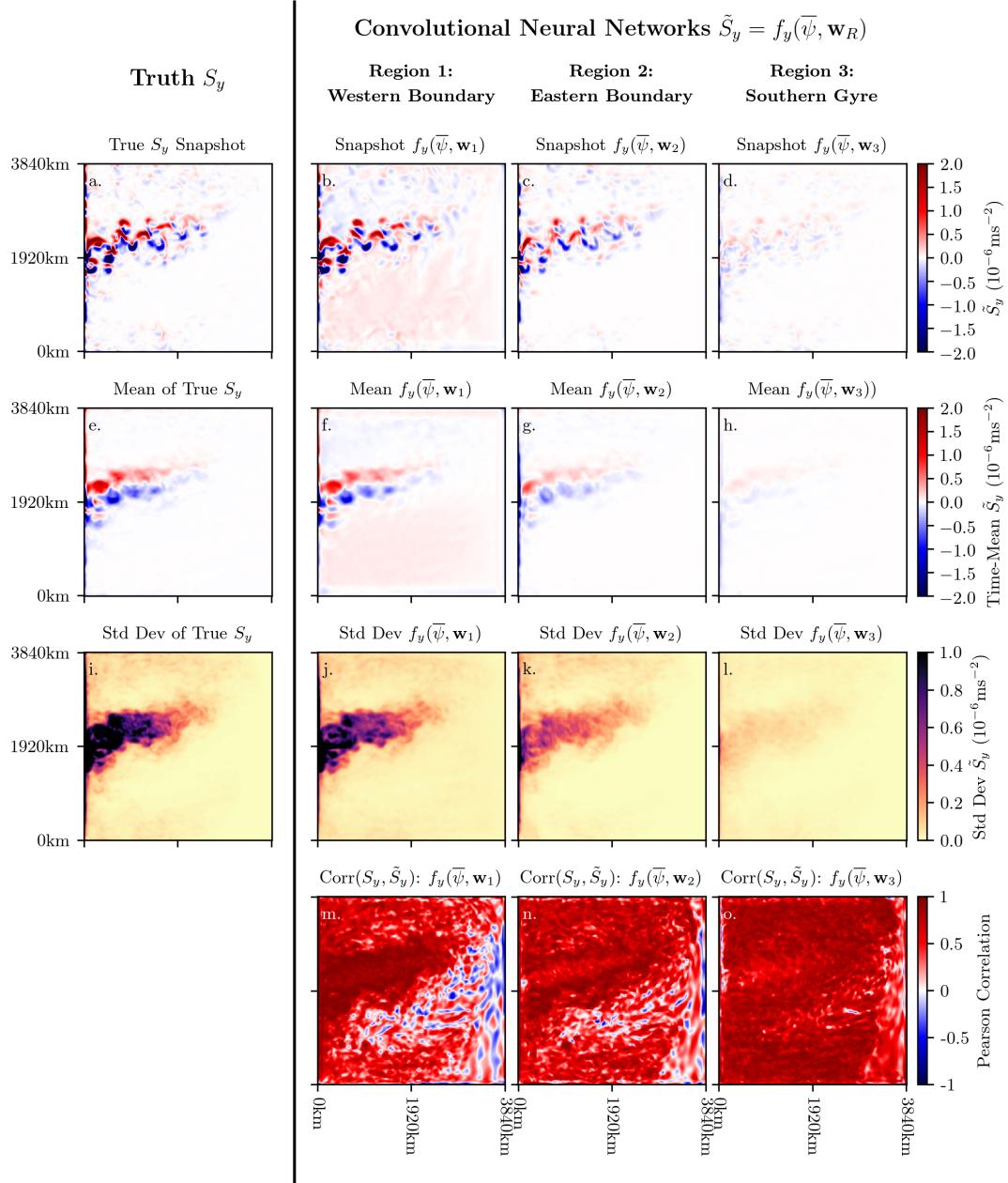
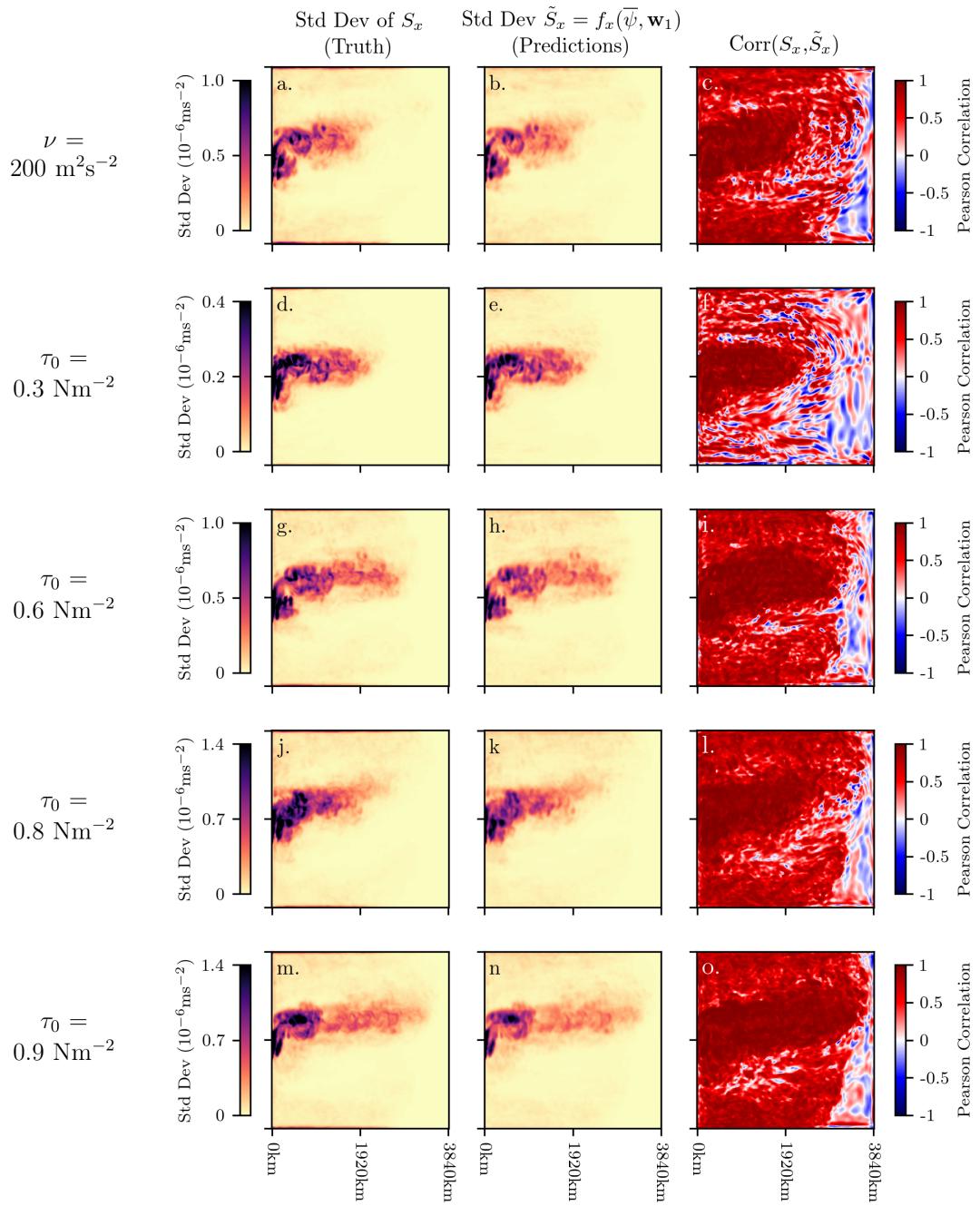
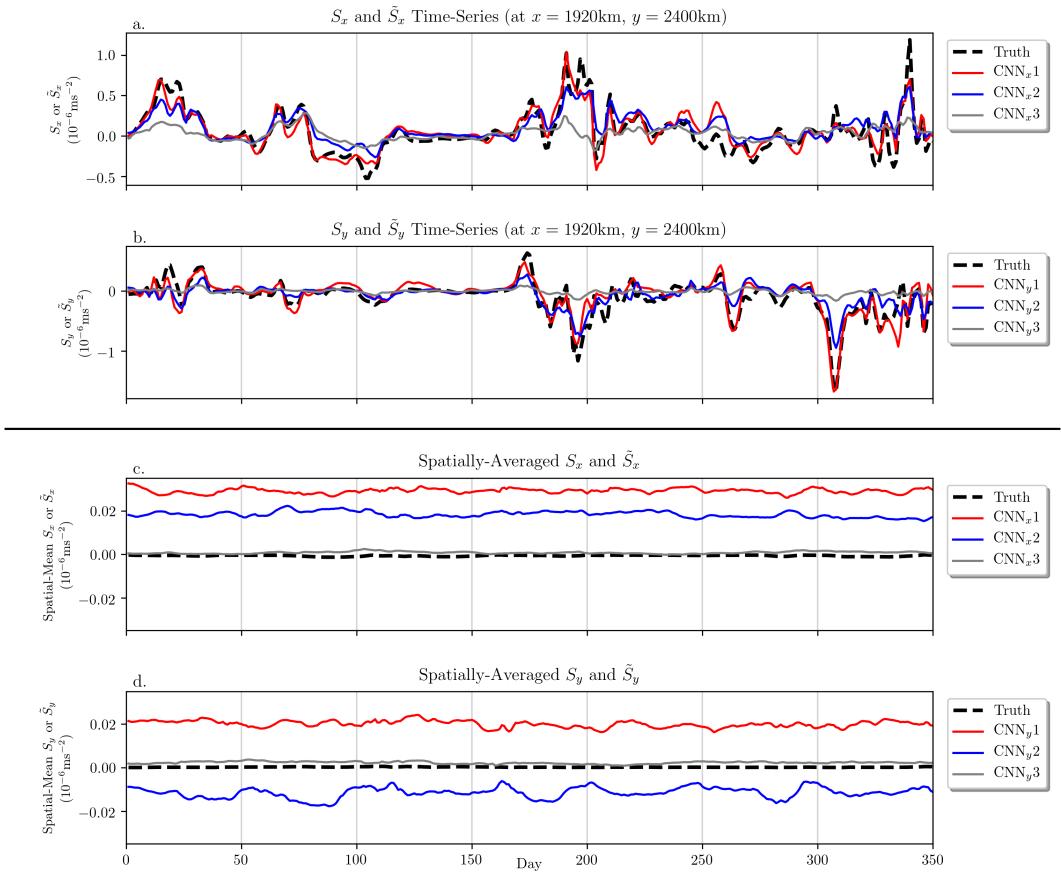


Figure 3. The same diagnostics as Figure 2, but for the meridional component of the sub-filter momentum forcing: the true S_y and the predictions \tilde{S}_y from the neural networks $f_y(\bar{\psi}, \mathbf{w}_1)$, $f_y(\bar{\psi}, \mathbf{w}_2)$, and $f_y(\bar{\psi}, \mathbf{w}_3)$.



792 **Figure 4.** Examining the ability to generalise to new regimes: using the trained neural network $f_x(\bar{\psi}, \mathbf{w}_1)$,
 793 we make predictions for model runs of different viscosities and wind forcings. From each model run, we
 794 use one year of the upper-layer filtered streamfunction to generate predictions \tilde{S}_x from $f_x(\bar{\psi}, \mathbf{w}_1)$ to see how
 795 they compare to the true S_x . We study a run of higher viscosity $\nu = 200 \text{ m}^2\text{s}^{-2}$, and runs with wind stress
 796 amplitude $\tau_0 = 0.3, 0.6, 0.8$, and 0.9 Nm^{-2} . Note that $f_x(\bar{\psi}, \mathbf{w}_1)$ was trained on a run with $\nu = 75 \text{ m}^2\text{s}^{-2}$
 797 and $\tau_0 = 0.8 \text{ Nm}^{-2}$, the standard deviation and correlation maps of which are included again here in Panels
 798 (j), (k), and (l).



799 **Figure 5.** Panels (a) and (b) show time series of the zonal and meridional components of the sub-filter
800 momentum forcing respectively, at a single point near the middle of the domain. Panels (c) and (d) also
801 show time series of the zonal and meridional components of the sub-filter momentum forcing, but this time
802 spatially-averaged over the entire domain.

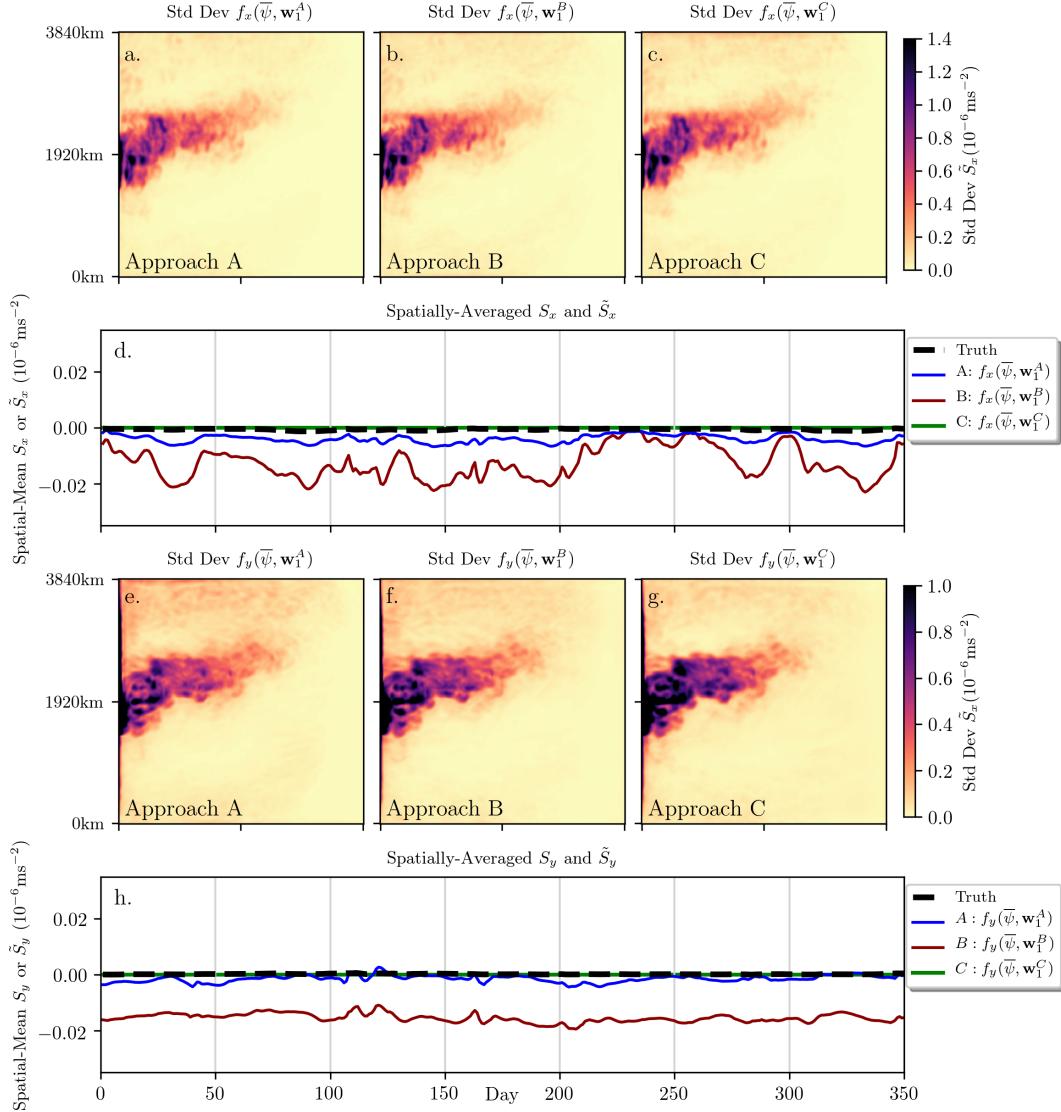
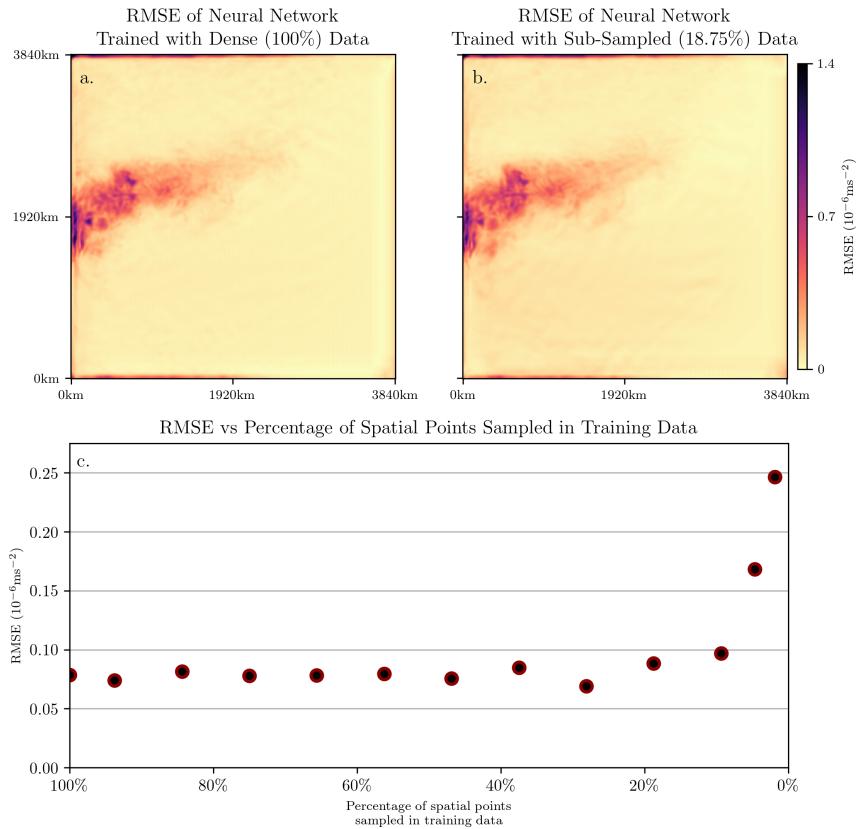


Figure 6. The standard deviation and spatial-average time series of the predictions \tilde{S}_x and \tilde{S}_y of the momentum conversing approaches A, B, and C. Panels (a), (b), and (c) show the standard deviation of \tilde{S}_x from $f_x(\bar{\psi}, \mathbf{w}_1^A)$, $f_x(\bar{\psi}, \mathbf{w}_1^B)$, 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).



808 **Figure 7.** Determining how under-sampling of the training data impacts neural network error. Panel (a)
809 shows the RMSE of the neural network $f_x(\bar{\psi}, \mathbf{w}_1)$ trained with dense (un-altered) training data, while
810 Panel (b) shows the RMSE of the neural network trained with sub-sampled (18.75%) data. Panel (c) shows the
811 RMSE as a function of the percentage of spatial points sampled at each time-slice of the training data. Note
812 that the RMSE is calculated over the full-domain during the validation period (the final year of data).

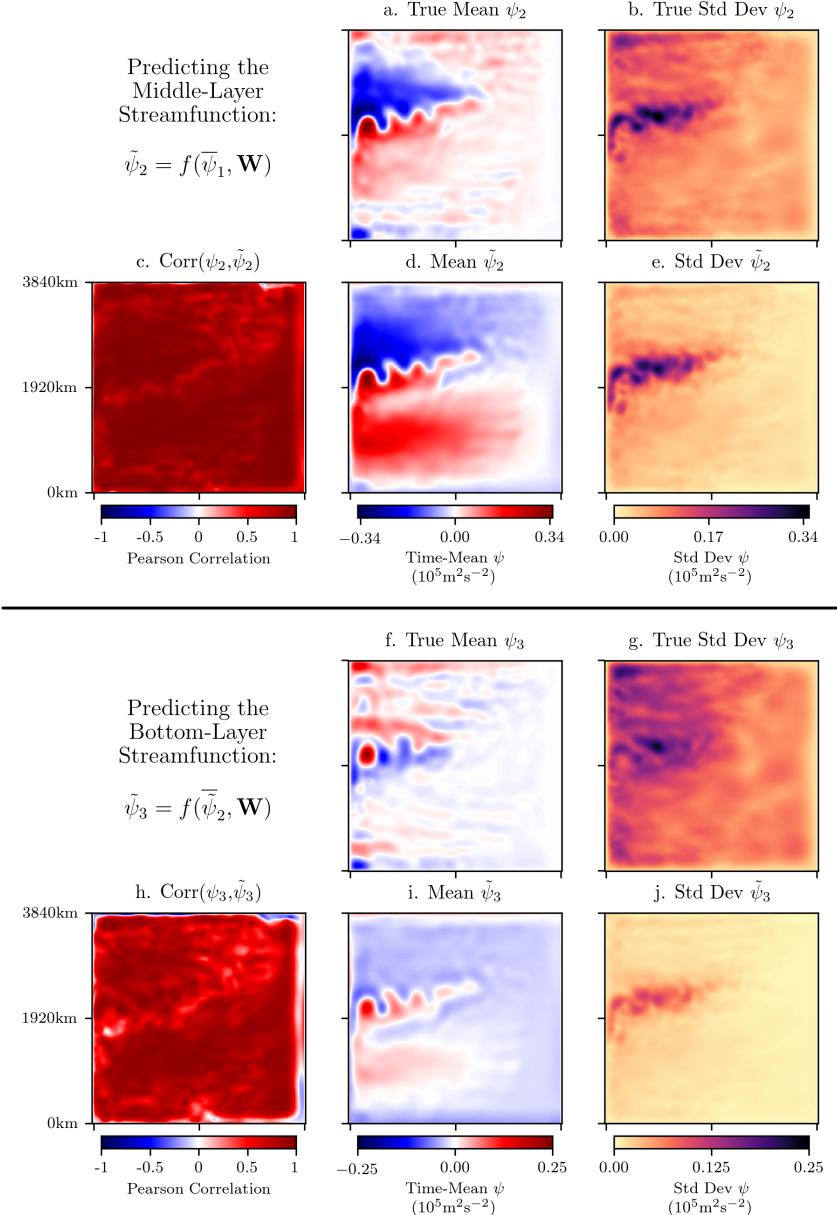


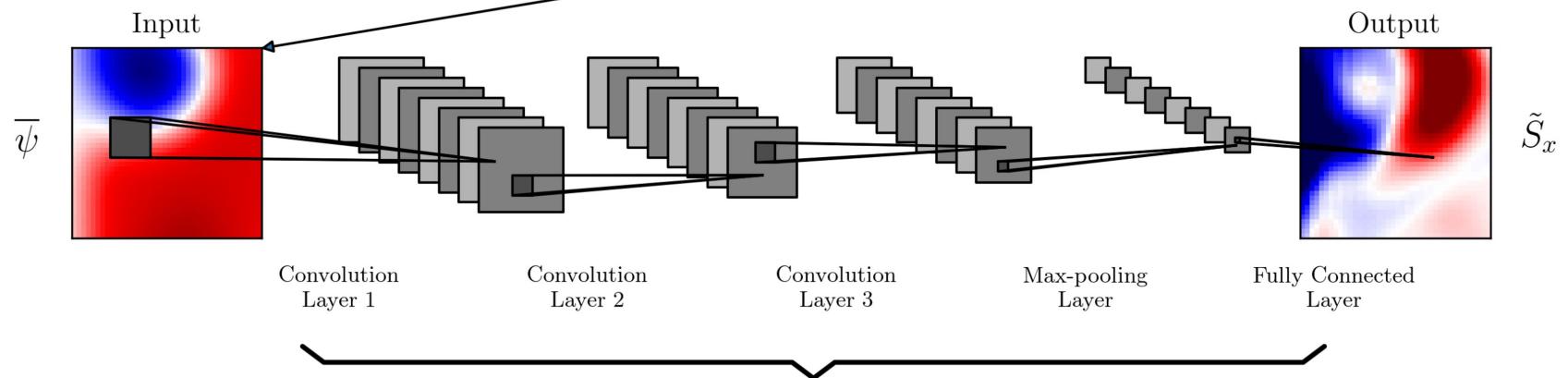
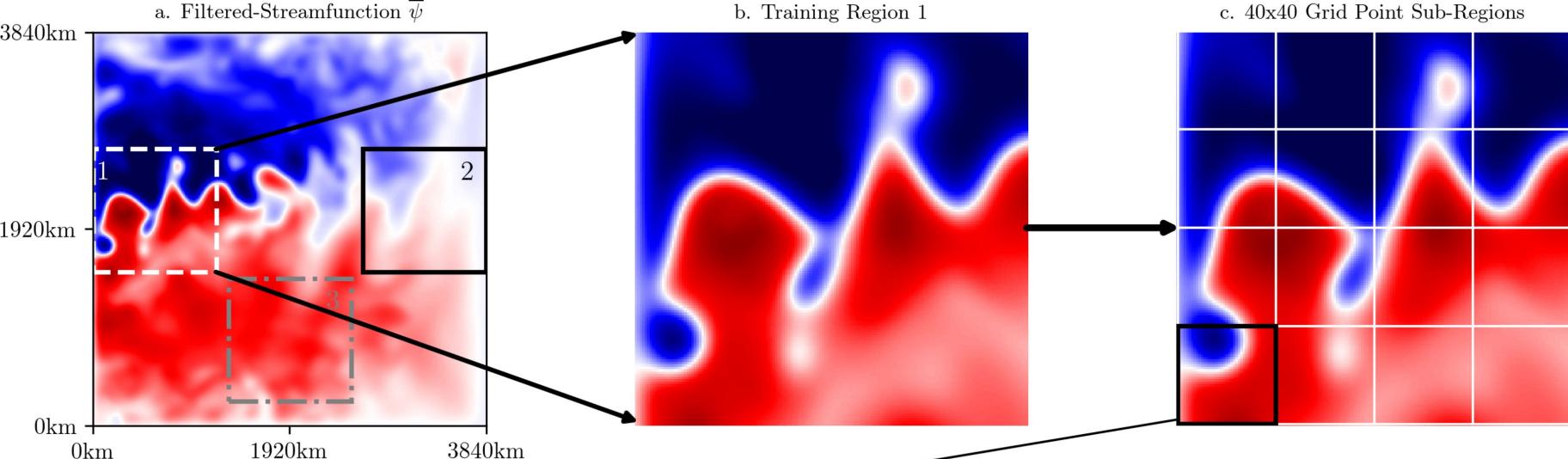
Figure 8. Predicting the middle- and bottom-layer streamfunctions ψ_2 and ψ_3 using the upper-layer filtered streamfunction $\bar{\psi}_1$. We first train a new neural network to predict ψ_2 from $\bar{\psi}_1$, i.e., $\psi_2 = f(\bar{\psi}_1, \mathbf{W})$; diagnostics of the true ψ_2 and the predictions $\tilde{\psi}_2$ are shown in the top-half of the Figure. We then take the same neural network that was trained to predict ψ_2 from $\bar{\psi}_1$, and now predict the bottom layer streamfunction ψ_3 using the predicted middle-layer streamfunction as the input, i.e., $\psi_3 = f(\tilde{\psi}_2, \mathbf{W})$; the diagnostics of the true ψ_3 and the predictions $\tilde{\psi}_3$ are shown in the bottom-half of the Figure.

- Arbic, B. K., K. L. Polzin, R. B. Scott, J. G. Richman, and J. F. Shriver (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., and C. S. Bretherton (2018), Prognostic validation of a neural network unified physics parameterization, *Geophysical Research Letters*.
- Chapman, C., and A. A. Charantonis (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., and J.-B. Sallée (2017), Can we reconstruct mean and eddy fluxes from argo floats?, *Ocean Modelling*, 120, 83–100.
- Chelton, D. B., M. G. Schlax, R. M. Samelson, and R. A. de Szoeke (2007), Global observations of large oceanic eddies, *Geophysical Research Letters*, 34(15).
- Chollet, F., et al. (2015), Keras, <https://keras.io>.
- Davis, R. E., M. D. Ohman, D. L. Rudnick, and J. T. Sherman (2008), Glider surveillance of physics and biology in the southern California current system, *Limnology and Oceanography*, 53(5part2), 2151–2168.
- Dong, C., C. C. Loy, K. He, and X. Tang (2016), Image super-resolution using deep convolutional networks, *IEEE transactions on pattern analysis and machine intelligence*, 38(2), 295–307.
- Durand, M., L.-L. Fu, D. P. Lettenmaier, D. E. Alsdorf, E. Rodriguez, and D. Esteban-Fernandez (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., G. de Souza Rolim, and A. S. Ferrando (2018), Rainfall prediction methodology with binary multilayer perceptron neural networks, *Climate Dynamics*, pp. 1–13.
- Gentine, P., M. Pritchard, S. Rasp, G. Reinaudi, and G. Yacalis (2018), Could machine learning break the convection parameterization deadlock?, *Geophysical Research Letters*.
- Giglio, D., V. Lyubchich, and M. Mazloff (2018), Estimating oxygen in the southern ocean using argo temperature and salinity, *Journal of Geophysical Research: Oceans*.
- Goodfellow, I., Y. Bengio, A. Courville, and Y. Bengio (2016), *Deep learning*, vol. 1, MIT press Cambridge.
- Greatbatch, R., X. Zhai, M. Claus, L. Czeschel, and W. Rath (2010a), Transport driven by eddy momentum fluxes in the Gulf Stream extension region, *Geophysical Research Letters*, 37(24).
- Greatbatch, R. J., X. Zhai, J.-D. Kohlmann, and L. Czeschel (2010b), 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., and J. Sun (2015), Convolutional neural networks at constrained time cost, in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5353–5360.
- Hewitt, H. T., M. J. Roberts, P. Hyder, T. Graham, J. Rae, S. E. Belcher, R. Bourdallé-Badie, D. Copsey, A. Coward, C. Guiavarch, 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., M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam (2017), Mobilenets: Efficient convolutional neural networks for mobile vision applications, *arXiv preprint arXiv:1704.04861*.

- 891 Jiang, G.-Q., J. Xu, and J. Wei (2018), A deep-learning algorithm of neural network for
 892 the parameterization of typhoon–ocean feedback in typhoon forecast models, *Geophysical*
 893 *Research Letters*.
- 894 Jochum, M., G. Danabasoglu, M. Holland, Y.-O. Kwon, and W. Large (2008), Ocean vis-
 895 cosity and climate, *Journal of Geophysical Research: Oceans*, 113(C6).
- 896 Kang, D., and E. N. Curchitser (2015), Energetics of eddy–mean flow interactions in the
 897 gulf stream region, *Journal of Physical Oceanography*, 45(4), 1103–1120.
- 898 Keating, S. R., and K. S. Smith (2015), Upper ocean flow statistics estimated from su-
 899 perresolved sea-surface temperature images, *Journal of Geophysical Research: Oceans*,
 900 120(2), 1197–1214.
- 901 Keating, S. R., A. J. Majda, and K. S. Smith (2012), New methods for estimating ocean
 902 eddy heat transport using satellite altimetry, *Monthly Weather Review*, 140(5), 1703–
 903 1722.
- 904 Kingma, D. P., and J. Ba (2014), Adam: A method for stochastic optimization, *arXiv*
 905 *preprint arXiv:1412.6980*.
- 906 Kjellsson, J., and L. Zanna (2017), The impact of horizontal resolution on energy transfers
 907 in global ocean models, *Fluids*, 2(3), 45.
- 908 Klambauer, G., T. Unterthiner, A. Mayr, and S. Hochreiter (2017), Self-normalizing neural
 909 networks, in *Advances in Neural Information Processing Systems*, pp. 972–981.
- 910 Krizhevsky, A., I. Sutskever, and G. E. Hinton (2012), Imagenet classification with deep
 911 convolutional neural networks, in *Advances in neural information processing systems*, pp.
 912 1097–1105.
- 913 Kutz, J. N. (2017), Deep learning in fluid dynamics, *Journal of Fluid Mechanics*, 814, 1–4.
- 914 Le Traon, P., and R. Morrow (2001), Ocean currents and eddies, in *International Geo-*
 915 *physics*, vol. 69, pp. 171–xi, Elsevier.
- 916 Le Traon, P., F. Nadal, and N. Ducet (1998), An improved mapping method of multisatel-
 917 lite altimeter data, *Journal of atmospheric and oceanic technology*, 15(2), 522–534.
- 918 LeCun, Y., Y. Bengio, and G. Hinton (2015), Deep learning, *nature*, 521(7553), 436.
- 919 Ling, J., R. Jones, and J. Templeton (2016a), Machine learning strategies for systems with
 920 invariance properties, *Journal of Computational Physics*, 318, 22–35.
- 921 Ling, J., A. Kurzawski, and J. Templeton (2016b), Reynolds averaged turbulence mod-
 922 elling using deep neural networks with embedded invariance, *Journal of Fluid Mechan-*
 923 *ics*, 807, 155–166.
- 924 Mana, P. P., and L. Zanna (2014), Toward a stochastic parameterization of ocean
 925 mesoscale eddies, *Ocean Modelling*, 79, 1–20.
- 926 McGovern, A., K. L. Elmore, D. J. Gagne, S. E. Haupt, C. D. Karstens, R. Lagerquist,
 927 T. Smith, and J. K. Williams (2017), Using artificial intelligence to improve real-time
 928 decision-making for high-impact weather, *Bulletin of the American Meteorological Soci-*
 929 *ety*, 98(10), 2073–2090.
- 930 Moeng, C.-H. (1984), A large-eddy-simulation model for the study of planetary boundary-
 931 layer turbulence, *Journal of the Atmospheric Sciences*, 41(13), 2052–2062.
- 932 Morrow, R., R. Coleman, J. Church, and D. Chelton (1994), Surface eddy momentum flux
 933 and velocity variances in the southern ocean from geosat altimetry, *Journal of Physical*
 934 *Oceanography*, 24(10), 2050–2071.
- 935 O’Gorman, P. A., and J. G. Dwyer (2018), Using machine learning to parameterize moist
 936 convection: potential for modeling of climate, climate change and extreme events, *Jour-*
 937 *nal of Advances in Modelling Earth Systems*.
- 938 Pathak, J., B. Hunt, M. Girvan, Z. Lu, and E. Ott (2018a), Model-free prediction of large
 939 spatiotemporally chaotic systems from data: a reservoir computing approach, *Physical*
 940 *review letters*, 120(2), 024,102.
- 941 Pathak, J., A. Wikner, R. Fussell, S. Chandra, B. R. Hunt, M. Girvan, and E. Ott (2018b),
 942 Hybrid forecasting of chaotic processes: using machine learning in conjunction with
 943 a knowledge-based model, *Chaos: An Interdisciplinary Journal of Nonlinear Science*,
 944 28(4), 041,101.

- Pope, S. (1975), A more general effective-viscosity hypothesis, *Journal of Fluid Mechanics*, 72(2), 331–340.
- Rocha, C. B., S. T. Gille, T. K. Chereskin, and D. Menemenlis (2016), Seasonality of submesoscale dynamics in the kuroshio extension, *Geophysical Research Letters*, 43(21), 11–304.
- Roemmich, D., G. C. Johnson, S. Riser, R. Davis, J. Gilson, W. B. Owens, S. L. Garzoli, C. Schmid, and M. Ignaszewski (2009), The argo program: Observing the global ocean with profiling floats, *Oceanography*, 22(2), 34–43.
- Rudnick, D. L., R. E. Davis, C. C. Eriksen, D. M. Fratantoni, and M. J. Perry (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*, Springer Science & Business Media.
- Scott, R. B., and F. Wang (2005), Direct evidence of an oceanic inverse kinetic energy cascade from satellite altimetry, *Journal of Physical Oceanography*, 35(9), 1650–1666.
- Simonyan, K., and A. Zisserman (2014), Very deep convolutional networks for large-scale image recognition, *arXiv preprint arXiv:1409.1556*.
- Su, H., W. Li, and X.-H. Yan (2018), Retrieving temperature anomaly in the global subsurface and deeper ocean from satellite observations, *Journal of Geophysical Research: Oceans*, 123(1), 399–410.
- Tracey, B. D., K. Duraisamy, and J. J. Alonso (2015), A machine learning strategy to assist turbulence model development, in *53rd AIAA Aerospace Sciences Meeting*, p. 1287.
- Vlachas, P. R., W. Byeon, Z. Y. Wan, T. P. Sapsis, and P. Koumoutsakos (2018), Data-driven forecasting of high-dimensional chaotic systems with long short-term memory networks, *Proc. R. Soc. A*, 474(2213), 20170,844.
- Waterman, S., and S. R. Jayne (2010), Eddy-mean flow interactions in the along-stream development of a western boundary current jet: An idealized model study, *Journal of Physical Oceanography*.
- Waterman, S., N. G. Hogg, and S. R. Jayne (2011), Eddy–mean flow interaction in the kuroshio extension region, *Journal of Physical Oceanography*, 41(6), 1182–1208.
- Zanna, L., P. P. Mana, J. Anstey, T. David, and T. Bolton (2017), Scale-aware deterministic and stochastic parametrizations of eddy-mean flow interaction, *Ocean Modelling*, 111, 66–80.
- Zanna, L., J. M. Brankart, M. Huber, S. Leroux, T. Penduff, and P. D. Williams (2018), Uncertainty and scale interactions in ocean ensembles: From seasonal forecasts to multidecadal climate predictions, *Quarterly Journal of the Royal Meteorological Society*.

Figure 1.



Neural network $\tilde{S}_x = f_x(\bar{\psi}, \mathbf{w}_1)$, trained to minimize loss $L \propto (S_x - \tilde{S}_x)^2$.

Figure 2.

Convolutional Neural Networks $\tilde{S}_x = f_x(\bar{\psi}, \mathbf{w}_R)$

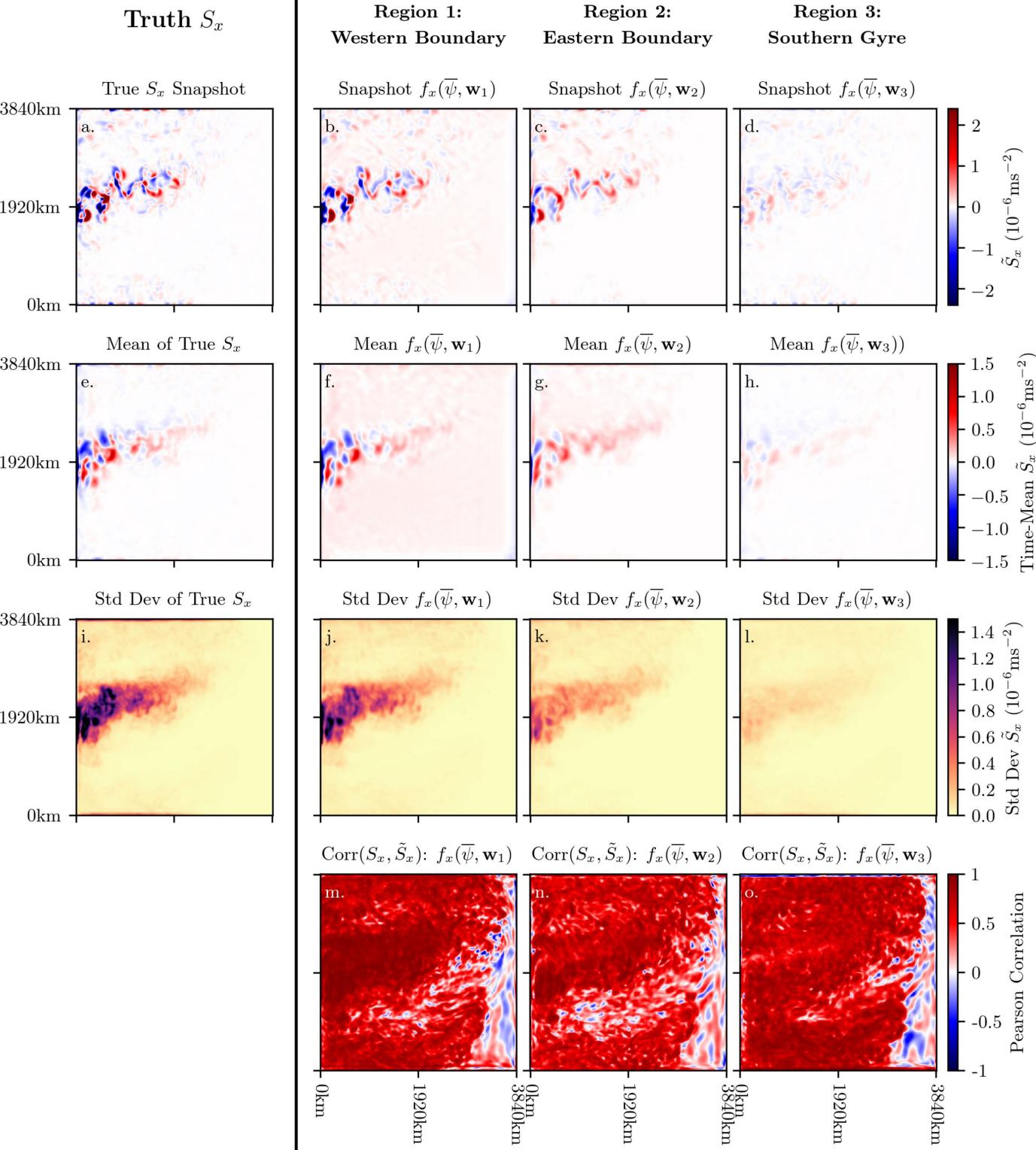


Figure 3.

Convolutional Neural Networks $\tilde{S}_y = f_y(\bar{\psi}, \mathbf{w}_R)$

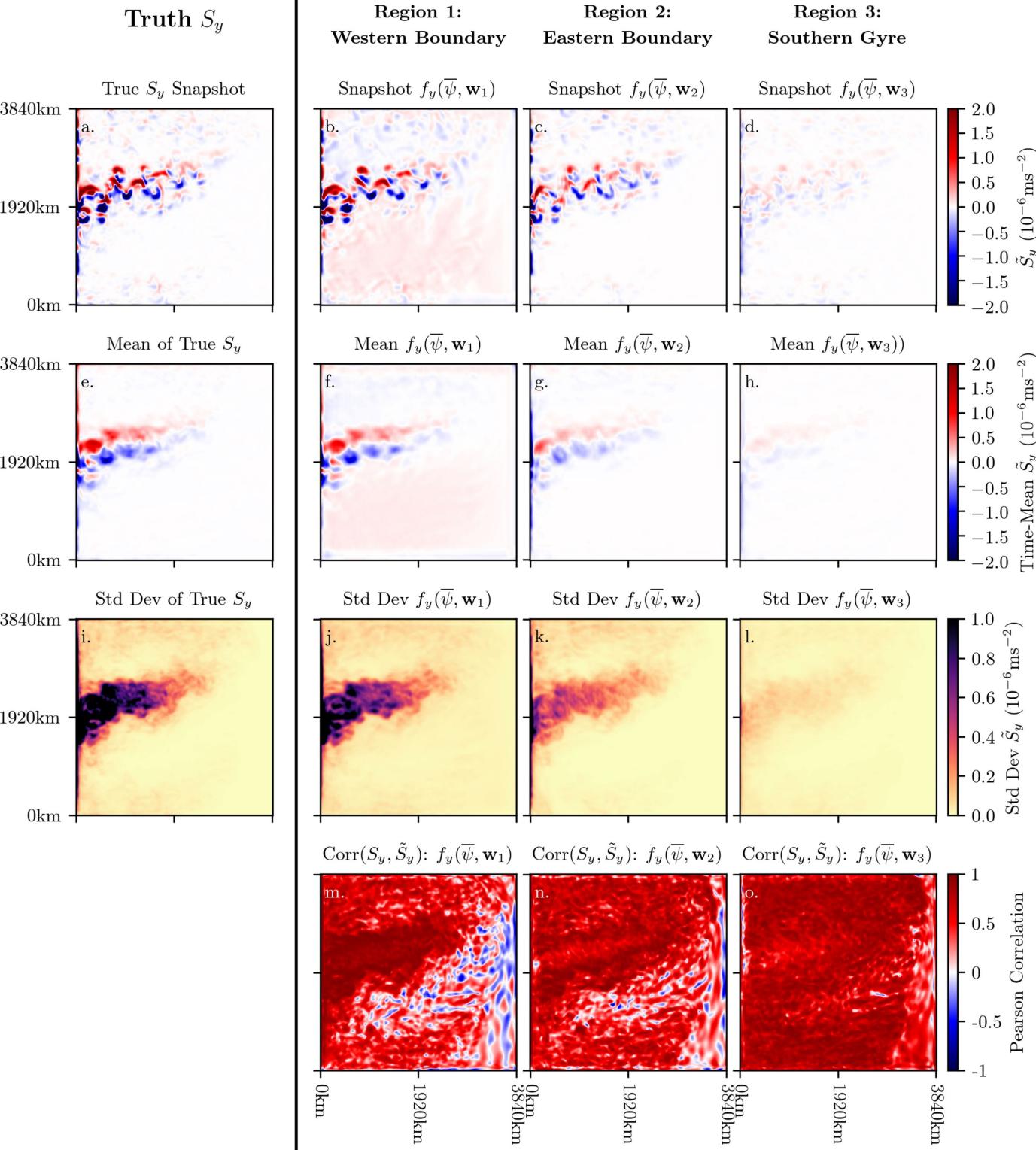


Figure 4.

$$\nu =$$

$$\tau_0 = \\ 0.3 \text{ Nm}^{-2}$$

$$\tau_0 =$$

$$\tau_0 =$$

$$\tau_0 =$$

Std Dev of S_x
(Truth)

Std Dev $\tilde{S}_x = f_x(\bar{\psi}, \mathbf{w}_1)$
(Predictions)

$$\text{Corr}(S_x, \tilde{S}_x)$$

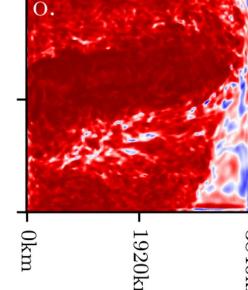
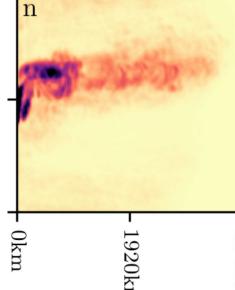
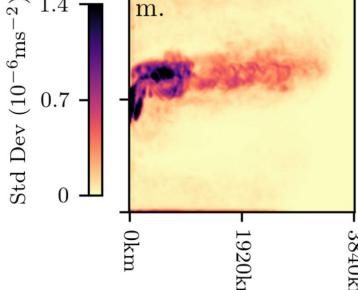
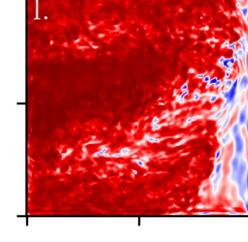
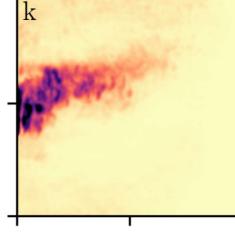
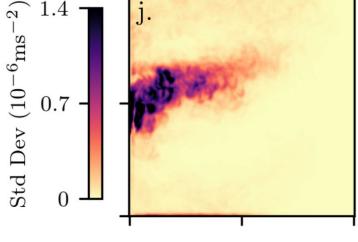
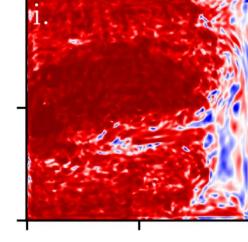
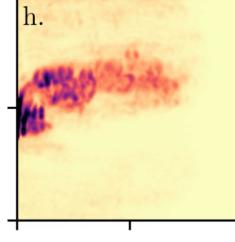
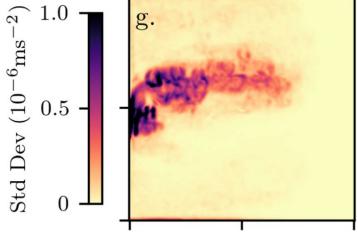
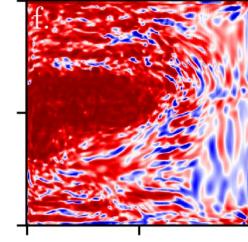
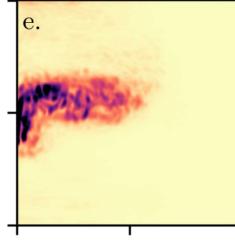
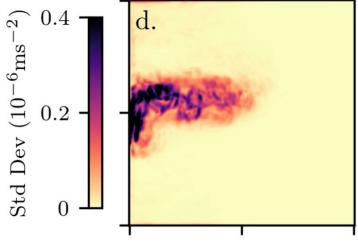
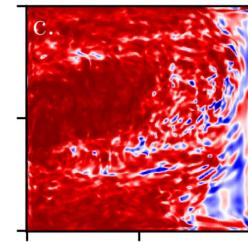
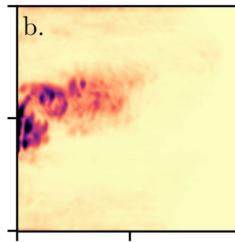
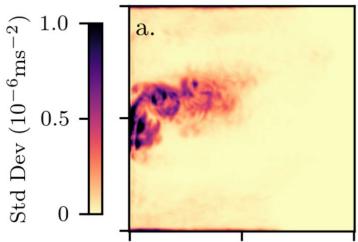
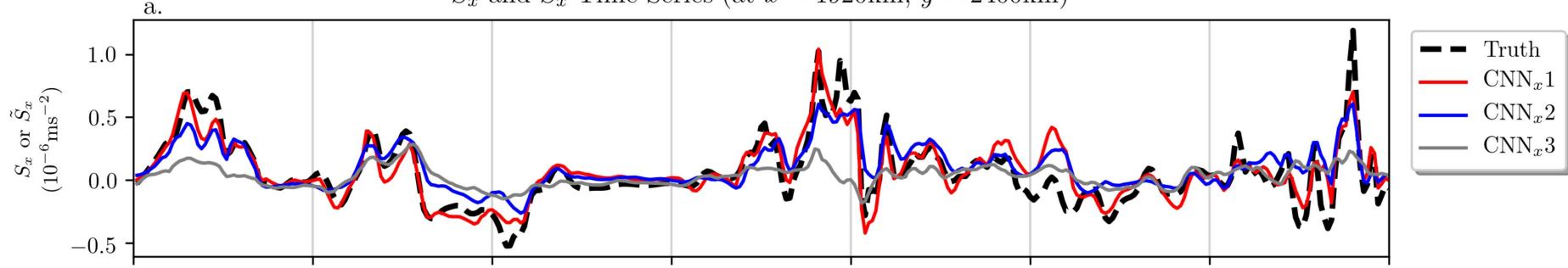
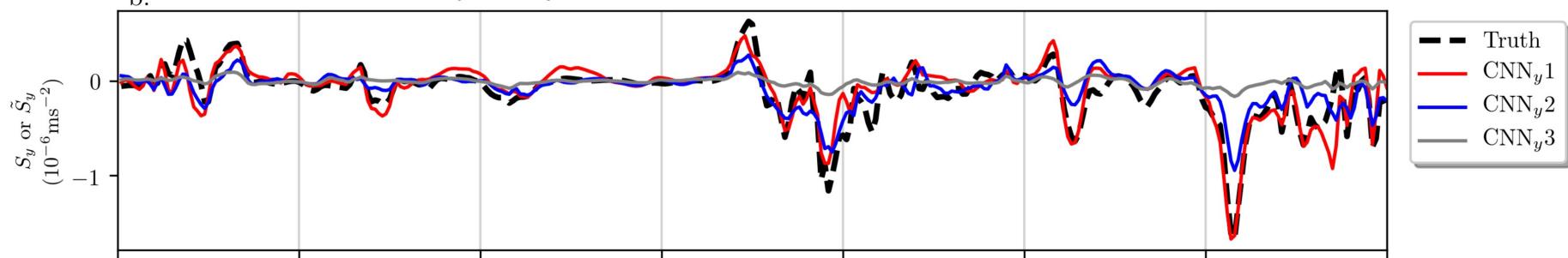


Figure 5.

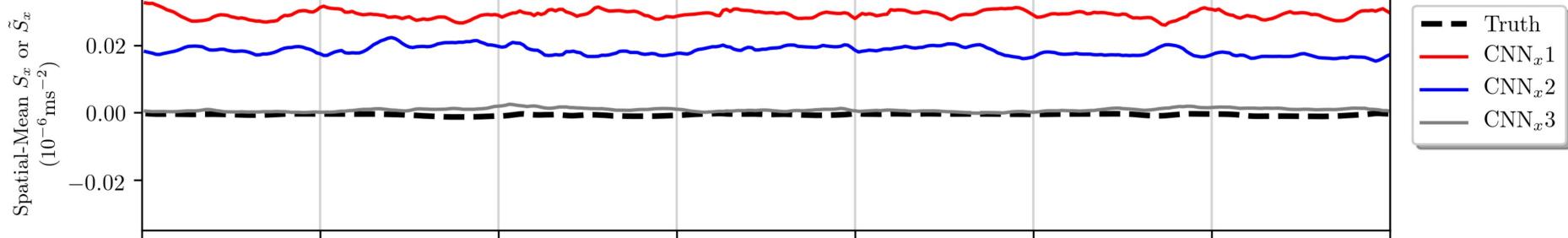
S_x and \tilde{S}_x Time-Series (at $x = 1920\text{km}$, $y = 2400\text{km}$)



S_y and \tilde{S}_y Time-Series (at $x = 1920\text{km}$, $y = 2400\text{km}$)



Spatial-Averaged S_x and \tilde{S}_x



Spatial-Averaged S_y and \tilde{S}_y

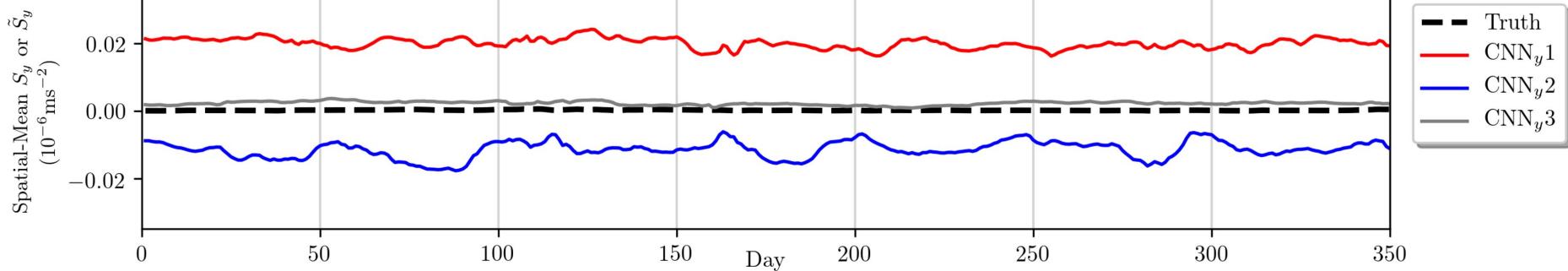


Figure 6.

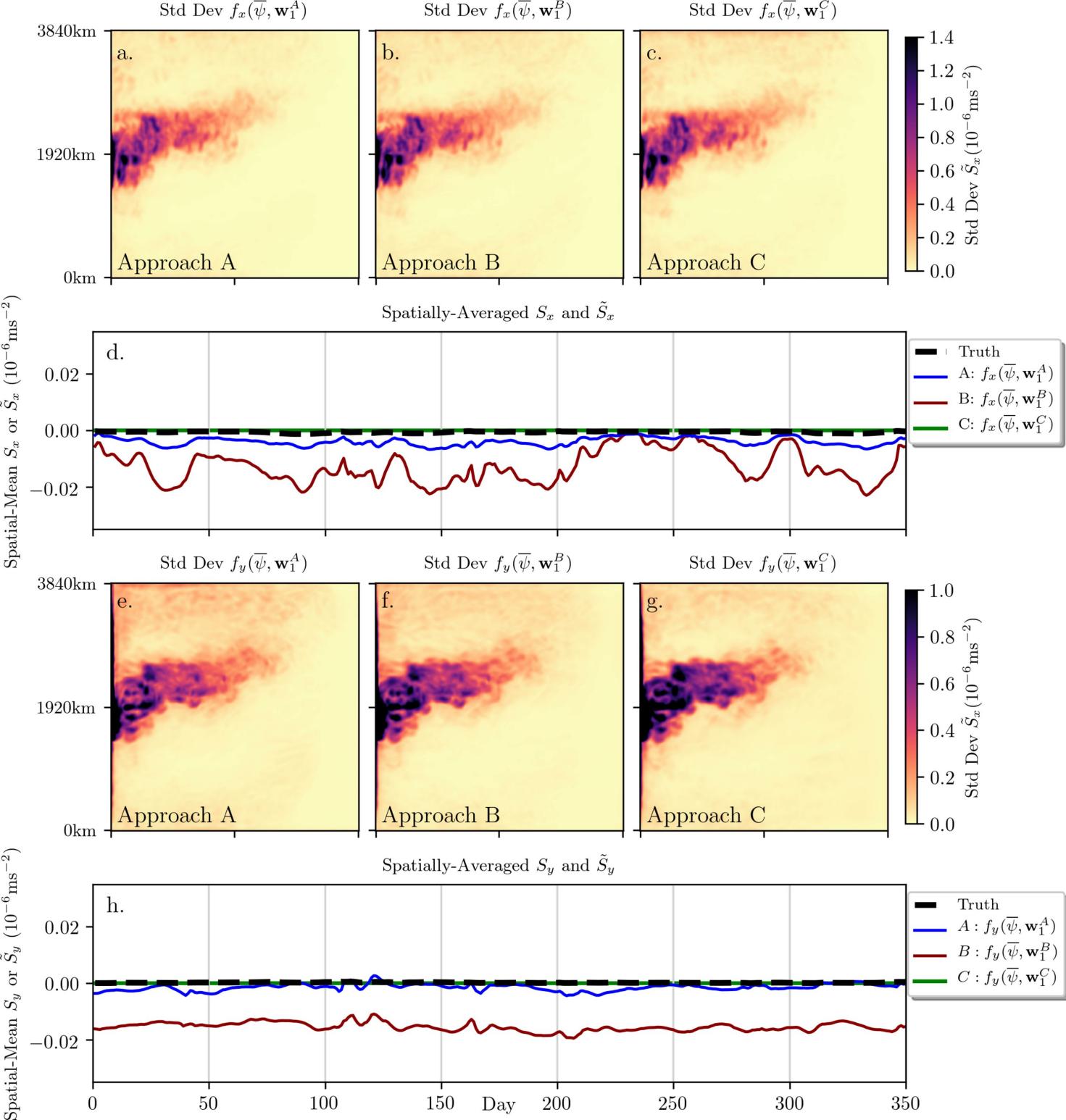
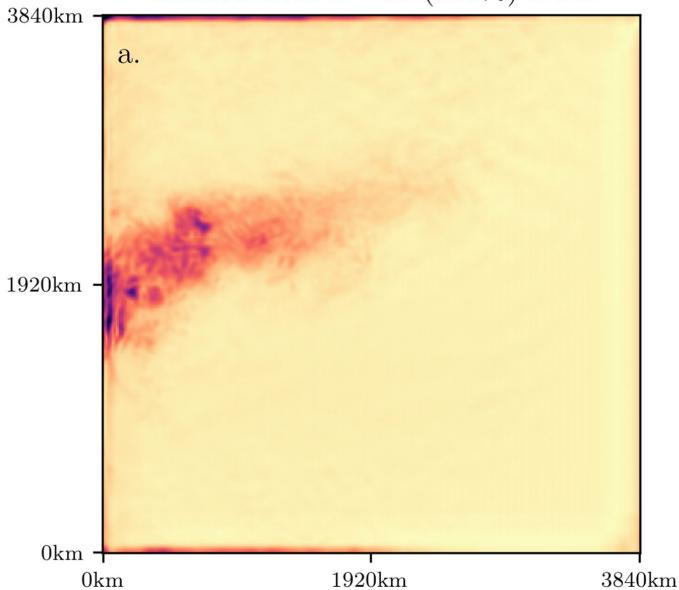
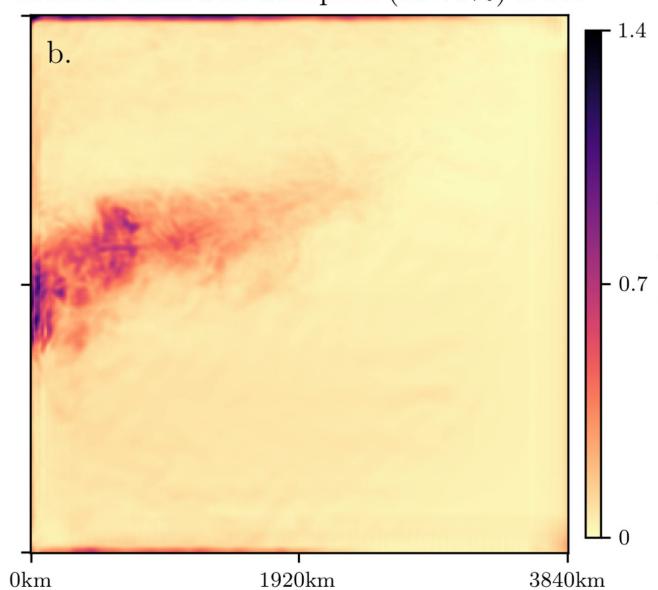


Figure 7.

RMSE of Neural Network
Trained with Dense (100%) Data



RMSE of Neural Network
Trained with Sub-Sampled (18.75%) Data



RMSE vs Percentage of Spatial Points Sampled in Training Data

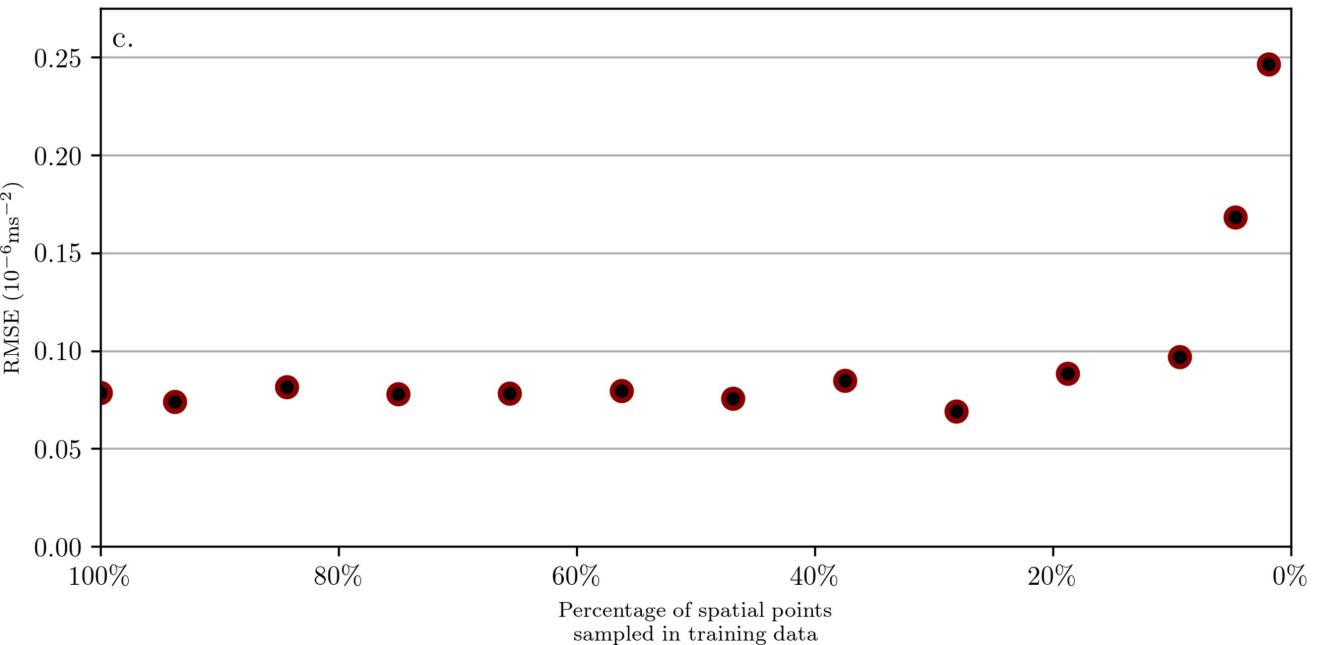
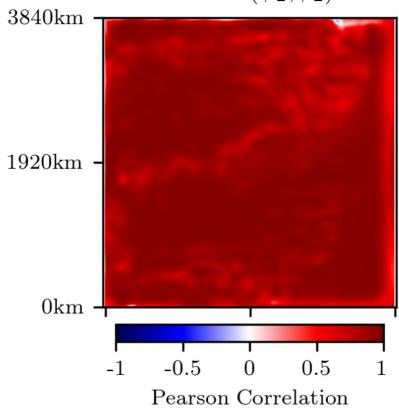


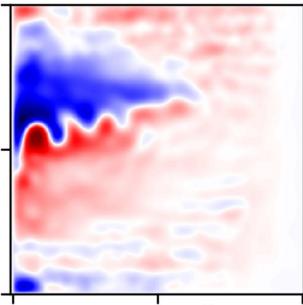
Figure 8.

Predicting the
Middle-Layer
Streamfunction:

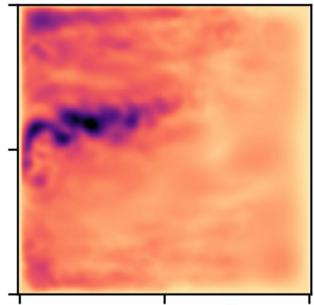
$$\tilde{\psi}_2 = f(\bar{\psi}_1, \mathbf{W})$$



a. True Mean ψ_2

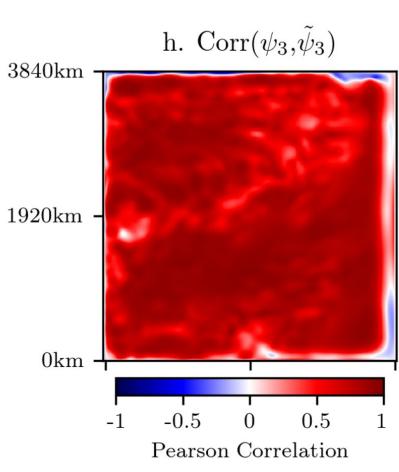


b. True Std Dev ψ_2

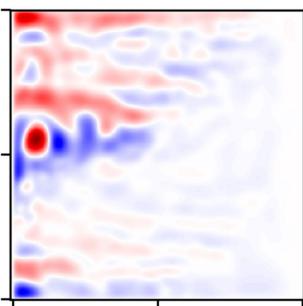


Predicting the
Bottom-Layer
Streamfunction:

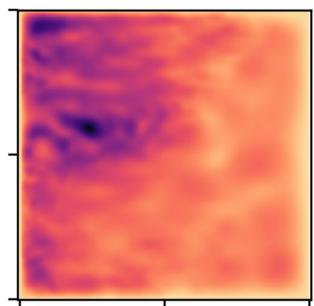
$$\tilde{\psi}_3 = f(\bar{\psi}_2, \mathbf{W})$$



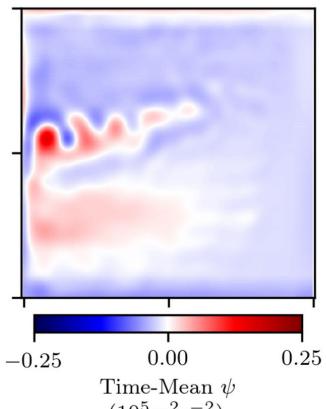
f. True Mean ψ_3



g. True Std Dev ψ_3



i. Mean $\tilde{\psi}_3$



j. Std Dev $\tilde{\psi}_3$

