

Estimating Ocean Surface Currents from Satellite Observable Quantities with Machine Learning

Anirban Sinha^{1,*}, Ryan Abernathey²

¹ California Institute of Technology, Pasadena, CA 91125, USA

² Lamont Doherty Earth Observatory of Columbia University, Palisades, NY, USA

Correspondence*:

Anirban Sinha
anirban@caltech.edu

2 ABSTRACT

Global surface currents are usually inferred from directly observed quantities like sea-surface height, wind stress by applying diagnostic balance relations (like geostrophy and Ekman flow), which provide a good approximation of the dynamics of slow, large-scale currents at large scales and low Rossby numbers. However, newer generation satellite altimeters (like the upcoming SWOT mission) will capture more of the high wavenumber variability associated with the unbalanced components, and applying these balances directly may lead to an incorrect unphysical estimate of the surface flow. In this study we explore Machine Learning (ML) algorithms as an alternate route to infer surface currents from satellite observable quantities. We train our ML models with SSH, SST and wind stress from available primitive equation ocean GCM simulation outputs as the inputs and make predictions of surface currents (u, v), which are then compared against the true GCM output. As a baseline example, we demonstrate that a linear regression model is ineffective at predicting velocities accurately beyond localized regions. In comparison, a relatively simple neural network (NN) can predict surface currents accurately over most of the global ocean, with lower mean squared errors than geostrophy+Ekman. Using a local stencil of neighboring grid points as additional input features, we can train the deep learning models to effectively “learn” spatial gradients and the physics of surface currents. By passing the stenciled variables through convolutional filters we can help the model learn spatial gradients much faster. Various training strategies are explored using systematic feature hold out and multiple combinations of point and stenciled input data fed through convolutional filters (2D / 3D), to understand the effect of each input feature on the NN’s ability to accurately represent surface flow. A model sensitivity analysis reveals that besides SSH, geographic information in some form is an essential ingredient required for making accurate predictions of surface currents with deep learning models.

Keywords: Neural Networks, Deep Learning, Surface Currents, Statistical models, Regression, geostrophy, Ekman flow

1 INTRODUCTION

The most reliable spatially continuous estimates of global surface currents in the ocean come from geostrophic balance applied to the sea surface height (SSH) field observed by satellite altimeters. For the most part, the dynamics of slow, large-scale currents (up to the mesoscale) are well approximated by

30 geostrophic balance, leading to a direct relationship between gradients of SSH and near-surface currents.
 31 However, current meter observations for the past few decades and some of the newer generation ultra-
 32 high-resolution numerical model simulations indicate the presence of an energized submesoscale as well
 33 as high-frequency waves / tides at smaller spatial and temporal scales (Rocha et al., 2016). These high-
 34 frequency unbalanced motions are likely to complicate the estimation of surface currents from from SSH
 35 in the upcoming Surface Water and Ocean Topography (SWOT) mission (Morrow et al., 2018). That is, the
 36 high-wavenumber SSH variability may represent a different, ageostrophic regime, where geostrophy might
 37 not be the best route to infer velocities. Motivated by this problem, in this study we explore statistical
 38 models based on machine learning (ML) algorithms for inferring surface currents from satellite observable
 39 quantities like SSH, wind and temperature. These algorithms offer a potential alternative to the traditional
 40 physics-based models.

41 The traditional method of calculating surface currents from sea surface height relies on the following
 42 physical principles. Assuming 2D flow and shallow water pressure, the momentum equation at the ocean
 43 surface can be written as:

$$\frac{\partial \mathbf{u}}{\partial t} + \mathbf{u} \cdot \nabla \mathbf{u} + f \times \mathbf{u} = -g \nabla \eta + \mathbf{F} \quad (1)$$

44 where \mathbf{F} is the frictional term due to wind stress. For a sufficiently low Rossby number (acceleration terms
 45 small; a questionable assumption for the submesoscale regime), the leading-order balances are geostrophy
 46 and Ekman flow. The surface flow can be split into a geostrophic and an ageostrophic, Ekman component
 47 ($\mathbf{u} = \mathbf{u}_g + \mathbf{u}_e$), and this leading-order force balance can be written as

$$f \times \mathbf{u}_g = -g \nabla \eta \quad (2)$$

$$f \times \mathbf{u}_e = F \quad (3)$$

48 Satellite altimetry products typically provide the sea surface height relative to the geoid (SSH, η), with
 49 tidally driven SSH signals removed (LeTraon and Morrow, 2001). The geostrophic velocities associated
 50 with the SSH anomalies are given by:

$$fv_g = g \frac{\partial \eta}{\partial x} \quad (4)$$

$$fu_g = -g \frac{\partial \eta}{\partial y} \quad (5)$$

51 Since geostrophic balance does not hold at the equator ($f \approx 0$), typically (Ducet et al., 2000), a higher
 52 order “equatorial geostrophic” treatment is used to compute velocities near the equator (Lagerloef et al.,
 53 1999), which is matched to the geostrophic regime away from the equator. Usually, the data-assimilative
 54 processing algorithms used to map along-track SSH observations to gridded maps (e.g. AVISO Ducet
 55 et al., 2000) also involve some form of temporal smoothing. The process of combining measurements from
 56 multiple satellites and filtering can also lead to spurious physical signals (Arbic et al., 2012) leading to
 57 exaggerated forward-cascades of energy.

58 In addition to the geostrophic velocities, some products like OSCAR (Ocean Surface Current Analysis
 59 Real Time, Bonjean and Lagerloef, 2002), or GEKCO (Geostrophic and Ekman Current Observatory,
 60 Sudre and Morrow, 2008; Sudre et al., 2013) provide an additional ageostrophic component due to Ekman
 61 flow. The Ekman velocity is related to friction, which in the upper layer of the ocean is provided by wind

62 stress and can be derived from the following equations:

$$fv_e + \frac{\partial \tau_x}{\partial z} = 0 \quad (6)$$

$$fu_e - \frac{\partial \tau_y}{\partial z} = 0 \quad (7)$$

$$\tau_x = \rho A_z \frac{\partial u}{\partial z} \quad (8)$$

$$\tau_y = \rho A_z \frac{\partial v}{\partial z} \quad (9)$$

63 Since the Coriolis parameter f changes sign at the equator, the functional relationship between velocity
64 and wind stress is different between the two hemispheres. In the Northern Hemisphere we derive:

$$u_e = \frac{1}{\rho \sqrt{2A_z|f|}}(\tau_x + \tau_y) \quad (10)$$

$$v_e = \frac{1}{\rho \sqrt{2A_z|f|}}(-\tau_x + \tau_y) \quad (11)$$

65 And in the Southern Hemisphere:

$$u_e = \frac{1}{\rho \sqrt{2A_z|f|}}(\tau_x - \tau_y) \quad (12)$$

$$v_e = \frac{1}{\rho \sqrt{2A_z|f|}}(\tau_x + \tau_y) \quad (13)$$

66 where A_z is the linear drag coefficient representing vertical eddy viscosity. Alternatively we can write these
67 equations in terms of the Ekman layer depth h_{Ek} which is related to the eddy viscosity A_z as:

$$h_{Ek} = \sqrt{\frac{2A_z}{f}} \quad (14)$$

68 Both of these quantities (A_z, h_{Ek}) are largely unknown for the global ocean and are estimated based on
69 empirical multiple linear regression from Lagrangian surface drifters (Lagerloef et al., 1999; Sudre et al.,
70 2013). Typical values of Ekman depth h_{Ek} in the ocean range from 10 to 40 meters .

71 So geostrophy + Ekman is the essential underlying physical/dynamical “model” currently used for
72 calculating surface currents from satellite observations. This procedure, combining observations with
73 physical principles, represents a top-down approach A more bottom-up approach would be a data driven
74 regression model that extracts information about empirical relationships from data. Recently, machine
75 learning (ML) methods have grown in popularity and have been proposed for a wide range of problems in
76 fluid dynamics: Reynolds-averaged turbulence models (Ling et al., 2016), detecting eddies from altimetric
77 SSH fields (Lguensat et al., 2017), reconstructing subsurface flow-fields in the ocean from surface fields
78 (Chapman and Charantonis, 2017; Bolton and Zanna, 2018), sub-gridscale modeling of PDEs (Bar-Sinai
79 et al., 2018), predicting the evolution of large spatio-temporally chaotic dynamical systems (Pathak et al.,
80 2018), parameterizing unresolved processes, like convective systems in climate models (Gentine et al.,
81 2018), or eddy momentum fluxes in ocean models (Bolton and Zanna, 2018), to name just a few examples.

82 In this study we aim to tackle a simpler problem than those cited above: training a ML model to “learn”
83 the empirical relationships between the different observable quantities (sea surface height, wind stress etc.)
84 and surface currents (u, v). The hypothesis to be tested is the following: Can we use machine learning to
85 provide surface current estimates that resolve small scale (balanced/unbalanced) turbulent processes better
86 than geostrophy+Ekman? The motivation for doing this exercise is two-fold:

- 87 1. It will help us understand how machine learning can be applied in the context of traditional physics-
88 based theories. ML is often criticised as a ”black box.” But can we use ML to complement our physical
89 understanding? This present problem serves as a good test-bed since the corresponding physical model
90 is straightforward and well understood.

91 2. It may be of practical value when SWOT mission launches.

92 We see this work as a stepping stone to more complex applications of ML to ocean remote sensing of ocean
93 surface currents.

94 This paper is organized as follows. In section 2, we introduce the dataset that was used, the framework
95 of the problem and identify the key variables that are required for training a statistical model to predict
96 surface currents. In section 3 we describe numerical evaluation procedure for baseline physics-based model
97 that we are hoping to match/beat. In sections 4 and 5 we discuss the statistical models that we used. We
98 start with the simplest statistical model - linear regression in Section 4 before moving on to more advanced
99 methods like neural networks in Section 5. In section 6 we summarize some the findings from the present
100 study, discuss some of the shortcomings of the present approach, propose some solutions as well as outline
101 some of the future goals for this project.

2 DATASET AND INPUT FEATURES

102 To focus on the physical problem of relating currents to surface quantities, rather than the observational
103 problems of spatio-temporal sampling and instrument noise, we choose to analyze a high-resolution global
104 general circulation model (GCM), which provides a fully sampled, noise-free realization of the ocean state.
105 The dataset used for this present study is the surface fields from the ocean component of the Community
106 Earth System Model (CESM), called the Parallel Ocean Program (POP) simulation (Smith et al., 2010)
107 which has a $\approx 0.1^\circ$ horizontal resolution, with daily-averaged outputs available for the surface fields. The
108 model employs a B-grid (scalars at cell centers, vectors at cell corners) for the horizontal discretization
109 and a three-time-level second-order-accurate modified leap-frog scheme for stepping forward in time. The
110 model solves the primitive equations of motion, which, for the surface flow, are essentially (1). Further
111 details about the model physics and simulations can be found in Small et al. (2014); Uchida et al. (2017).

112 We selected this study because of the long time record of available data (approx. 40 years), although, in
113 retrospect, we found that all our models can be trained completely with just a few days of output!

114 A key choice in any ML application is the choice of features, or inputs, to the model. In this paper,
115 we experiment with a range of different feature combinations; seeing which features are most useful for
116 estimating currents is indeed one of our aims. The features we choose are all quantities that are observable
117 from satellites: SSH, surface wind stress (τ_x and τ_y), sea-surface temperature (SST, θ) and sea-surface
118 Salinity (SSS). Our choice of features is also motivated by the traditional physics-based model: the same
119 information that goes into the physics-based model should also prove useful to the ML model Just like the
120 physics-based model, all the ML models we consider are pointwise, local models: the goal is to predict the
121 2D velocity vector u, v at each point, using data from at or around that point.

122 Beyond these observable physical quantities, we also need to provide the models with geographic
 123 information about the location and spacing between the neighboring points. In the physics-based model,
 124 geography enters in two places: 1) in the Coriolis parameter f , and 2) in the grid spacing dx and dy , which
 125 varies over the model domain. Geographic information can be provided to the statistical models in a few
 126 different ways. The first method involves providing the same kind of spatial information that is provided to
 127 the physical models, *i.e.* f and local grid spacings - dx and dy . We can also encode geographic information
 128 (lat , lon) in our input features, using a coordinate transformation of the form:

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} \sin(\text{lat}) \\ \sin(\text{lon}) \cdot \cos(\text{lat}) \\ -\cos(\text{lon}) \cdot \cos(\text{lat}) \end{bmatrix} \quad (15)$$

129 to transform the spherical polar lat-lon coordinate into a homogeneous three dimensional coordinate
 130 (Gregor et al., 2017). This transformation gives the 3D position of each point in Euclidean space, rather
 131 than the geometrically warped lat / lon space (which has a singularity at the poles and a discontinuity at the
 132 dateline). Note that one of the coordinates - X , that comes out of this kind of coordinate transformation, is
 133 functionally the same as the coriolis parameter (f) normalized by 2Ω (Ω = Earth's rotation). Therefore
 134 we will use X as proxy for f for all the statistical models throughout this study. We also explored another
 135 approach where the only geographic information provided to the models is X ($= \frac{f}{2\Omega}$).

136 Since geostrophic balance involves spatial derivatives, it is not sufficient to simply provide SSH and the
 137 local coordinates pointwise. In order to compute derivatives, we also need the SSH of the surrounding
 138 grid points as a local stencil around each grid point. The approach we used for providing this local stencil
 139 is motivated by the horizontal discretization of the POP model. Horizontal derivatives of scalars (like
 140 SSH) on the B grid requires 4 cell centers. At every timestep, each variable of the The 1° POP model
 141 ouput has 3600×2400 data points (minus the land mask). We can simply rearrange each variable as a
 142 $1800 \times 1200 \times 2 \times 2$ dataset or split it into 4 variables each with 1800×1200 data points, corresponding
 143 to the 4 grid cells required for taking spatial derivatives. The variables that require spatial a spatial stencil
 144 for physical models, we will refer to as the stencil inputs. For the variables for which we do not need
 145 spatial derivatives for (like wind stress), we can simply use every alternate grid point resulting in a dataset
 146 of size 1800×200 . We will refer to these variables as point inputs. For the purpose of the statistical
 147 models the inputs need to be flattened and have all the land points removed. This means that each input
 148 variable has a shape of either $N \times 2 \times 2$ or N depending on whether or not a spatial stencil is used (where
 149 $N = 1800 \times 1200$ – the points that fall over land). Alternatively we can think of the stencilled variable
 150 as 4 features of length N . This kind of stencil essentially coarsens the resolution of the targets, and point
 151 variables.

152 Similarly we can also construct a 3 point time stencil, by providing the values at preceding and succeeding
 153 time steps as additional inputs so that each variable that is stencilled in space and time has a shape of
 154 $N \times 2 \times 2 \times 3$ (or 12 features of length N).

155 This data preparation leads to 10 potential features (for some of which we will use a stencil, which further
 156 expands the feature vector space) for predicting u, v at each point : τ_x, τ_y , SSH (η), SST (θ), SSS (S), the 3
 157 transformed coordinates (X, Y, Z) and the local grid spacings (dx and dy).

158 For building any statistical / ML model, we need to split the dataset into 2 main parts, *i.e.* training and
 159 testing. For the purpose of training our machine learning models, the first step involves extracting the above
 160 mentioned variables from the GCM output as the input features and the GCM output surface velocities

161 u, v as targets for the ML model. The data extracted from the GCM output for a certain date (or range of
 162 dates) is then used to fit the model parameters. This part of the dataset is called the training dataset. During
 163 training, the model minimizes a chosen cost function (we used mean absolute error for our experiments,
 164 but using mean squared error produced very similar results) and typically involves a few passes through
 165 this section of dataset. The trained models are then used to make predictions of u, v for a different date (or
 166 range of dates) where the model only receives the input variables. The model predictions are evaluated by
 167 comparing with the true (GCM output) velocity fields for that particular date (date range). This part of the
 168 dataset, which the model has not seen during training, that is used to evaluate model predictions is called
 169 the test dataset.

3 BASELINE PHYSICS-BASED MODEL: GEOSTROPHY + EKMAN

170 The two components of the physics-based model used as the baseline for our ML models are geostrophy
 171 and Ekman flow. In this section we describe how these two components are numerically evaluated for our
 172 dataset. For the sake of fair comparison, we evaluate the geostrophic and Ekman velocities from the same
 173 features that are provided to the regression models. With the POP model's horizontal discretization,
 174 finite-difference horizontal derivatives and averages are defined as (Smith et al., 2010) :

$$\psi_x = [\psi(x + \Delta_x/2) - \psi(x - \Delta_x/2)] / \Delta_x \quad (16)$$

$$\overline{\psi^x} = [\psi(x + \Delta_x/2) + \psi(x - \Delta_x/2)] / 2 \quad (17)$$

175 With the data preparation and stencil approach described in the previous section, η now has a shape of
 176 $N \times 2 \times 2$ and the f, u, v, dx, dy are all variables of length N . Following (4) and (5) the geostrophic
 177 velocities (u_g^j, v_g^j) are calculated on the stencil as as:

$$v_g^j = g/f^j [\eta_i(1, 1) + \eta_i(0, 1) - \eta_i(1, 0) - \eta_i(0, 0)] / 4dx^j \quad (18)$$

$$u_g^j = -g/f^j [\eta^j(1, 1) + \eta^j(1, 0) - \eta^j(0, 1) - \eta^j(0, 0)] / 4dy^j \quad (19)$$

$$(20)$$

178 where $j \in [1, N]$. Similarly the Ekman velocity is calculated numerically from the τ_x^j, τ_y^j and f^j as

$$u_e^j = \begin{cases} \frac{1}{\rho\sqrt{2A_z|f^j|}}(\tau_x^j + \tau_y^j), & \text{if } f^j > 0 \\ \frac{1}{\rho\sqrt{2A_z|f^j|}}(\tau_x^j - \tau_y^j), & \text{if } f^j < 0 \end{cases} \quad v_e^j = \begin{cases} \frac{1}{\rho\sqrt{2A_z|f^j|}}(-\tau_x^j + \tau_y^j), & \text{if } f^j > 0 \\ \frac{1}{\rho\sqrt{2A_z|f^j|}}(\tau_x^j + \tau_y^j), & \text{if } f^j < 0 \end{cases} \quad (21)$$

179 For calculating the Ekman velocity, we used constant values for vertical diffusivity ($A_z = 8 \times 10^{-3} m^2/s$)
 180 and density of water at the surface, ($\rho = 1027 kg/m^3$). It should be noted that both these quantities vary
 181 both spatially and temporally in the real ocean. For the vertical diffusivity we came up with this estimate
 182 by solving for A_z that provides the best fit between zonal mean $((u, v)_{true} - (u, v)_g)$ and $(u, v)_e$. In the
 183 CESM high res POP simulations, the parameterized vertical diffusivity was capped around $100 cm^2/s$
 184 (Smith et al., 2010). For plotting spatial maps for both the physics based model predictions as well as the
 185 statistical model predictions, the velocity fields are then reshaped into 1800×1200 arrays, after inserting
 186 the appropriate land masks.

4 MULTIPLE LINEAR REGRESSION MODEL

187 The simplest of all statistical prediction models is essentially multiple linear regression, where an output or
 188 target is represented as some linear combination of the inputs. The input is characterized by a feature vector
 189 \mathbf{x}_i^j where $i \in [1, n_f]$; $j \in [1, N]$, N being the number of samples, and n_f being the number of features.
 190 We can now write the linear regression problem as $\mathbf{U}^j = \mathbf{x}_i^j \cdot \beta_i + \delta^j$. where β_i are the coefficients or
 191 weight vector. For our regression problem, the input features are wind stress, sea surface height and the
 192 3 dimensional transformed coordinates. Of those features, η , X , Y , Z are stencil inputs (meaning 4 input
 193 columns per feature) and τ_x , τ_y are the point inputs, resulting in a total of 18 input features. The aim
 194 therefore is to find the coefficients β_i that minimize the loss (error) represented by δ^j for a training set of \mathbf{x}_i^j
 195 and \mathbf{U}^j (\mathbf{x}_{itrain}^j , \mathbf{U}_{train}) and use these coefficients for a test set of \mathbf{x}_i^j (\mathbf{x}_{itest}^j) to make predictions for \mathbf{U}^j
 196 (\mathbf{U}_{pred}^j). For implementing linear regression model as well as the deep learning models discussed in this
 197 study we use the Python library Keras (<https://keras.io>) (Chollet et al., 2015), a high-level wrapper around
 198 TensorFlow (<http://www.tensorflow.org>).

199 Linear regression can be performed in one of 2 different ways

- 200 • The matrix method or Normal equation method (where we solve for the coefficients β that minimize
 201 the squared error $\|\delta\|^2 = \|\mathbf{U} - \mathbf{X}^T \cdot \beta\|^2$ and involves computing the pseudo-inverse of $\mathbf{X}^T \cdot \mathbf{X}$).
- 202 • A stochastic gradient descent (SGD) method (which represents a more general procedure that can be
 203 used for different regression algorithms with different choices for optimizers and is more scalable for
 204 larger datasets).

205 The normal-equation method is less computationally tractable for large datasets (large number of samples)
 206 since it requires loading the full dataset into memory for calculating the pseudoinverse of $\mathbf{x}_i^j \cdot \mathbf{x}_i^j$, whereas
 207 the SGD method works well even for large datasets, but requires tuning of the learning rate. Due to the
 208 versatility offered by the gradient descent method we used that for performing the linear regression although
 209 the normal equation method also produced similar results. The essential goal for any regression problem is
 210 to minimize a predetermined cost / loss function (which for our experiments we chose as the mean absolute
 211 error) :

$$J = MAE = \overline{(|u_{pred} - u_{true}| + |v_{pred} - v_{true}|)} \quad (22)$$

212 where the overbar denotes the average over all samples. Fig. 2(a) shows a schematic of the linear regression
 213 model. The number of trainable parameters for our example with 18 inputs and 2 outputs is 38 (18×2
 214 weights + 2 biases). For this as well as all subsequent models discussed here, we used the same optimizer
 215 (Adam, (Kingma and Ba, 2017)) and loss function (Mean absolute error, MAE). All models are trained on
 216 1 day of GCM output data and we use the same date of model output as the training data for all models.

217 We start by splitting the global ocean into 3 boxes to zoom into three distinct regions of dynamical
 218 importance in oceanography, namely the Gulf stream, Kuroshio and Southern ocean / Antarctic circumpolar
 219 current (ACC). The Kuroshio region is chosen to extend south of the equator to include the equatorial jets
 220 as well as to test whether the models can generalize to large variations in f. The daily averaged GCM output
 221 surface speed on a particular reference date, with the three regions (marked by three different colored
 222 boxes) is shown in Fig. 1. We then train 3 different linear regression models with training data from these
 223 three sub-domains. We also trained a linear regression model for the whole globe using the same model
 224 architecture. During training, the models are fed a shuffled batch of the training data with 32 samples in
 225 each batch and the loss (MAE) is computed for the batch. For the linear regression model as well as for
 226 all the neural networks discussed in this study we present here we kept the batch size constant. Changing

227 the batch size does not significantly alter the loss at the end of training, but smaller batch sizes generally
228 help the model learn faster. The different models, the number of epochs (an epoch is defined as one pass
229 through the training dataset) used for each, and losses at the end of training and during evaluation against
230 a test dataset are summarized in table 1. The evolution of model loss function during training for the 3
231 different models are presented in Fig. 3. Linear regression is shown in the darker colors. The big jumps
232 in the loss function correspond to the end of an epoch. We plot the models' training progress in the Gulf
233 Stream region for 8 epochs, and for 5 epochs on the Kuroshio and ACC regions. The trained models are
234 then evaluated for a test dataset (which the model has not seen, GCM output from a different point in time)
235 and the evaluation loss is plotted as the horizontal dashed lines. The linear regression model trained on
236 the whole globe is also evaluated for each subdomain (gulf stream, Kuroshio, ACC) and the global model
237 evaluation loss is plotted as the dotted line. Comparing the model losses in the 3 different sections, we find
238 that the linear regression model performs the most poorly for the Kuroshio region (i.e the subdomain with
239 the most variation in f). The model does progressively better for the gulf stream and the ACC in terms
240 of MAE, where the variations in f are relatively smaller in comparison. However, the root mean squared
241 error of predicted velocities is still quite large in all these regions (second panels of Figs. 4, 5, 6). The
242 linear regression model trained on the global ocean does even worse during evaluation. Since geostrophy
243 relies on non-linear combination of the Coriolis parameter (f) with the spatial gradients, linear regression
244 is ineffective at predicting velocities beyond localized regions with small variation of f or little mesoscale
245 activity. This shows that a linear model fails to accurately represent surface currents in any region that
246 includes significant variation in the Coriolis parameter f . Even in regions far enough from the equator
247 such that the variation in f is not significant (like the gulf stream or ACC), the performance of such a
248 linear model does not improve with more training examples and/or starts overfitting. We also show that a
249 lower MAE during training does not necessarily guarantee that the model is picking up on the small scale
250 fluctuations in velocity, as can be seen from the relatively large squared errors especially in and around
251 high surface current regions (Figs. 4, 5, 6). We suspect that this failure is largely due to the fact that the
252 linear model is trying to fit the velocities as a linear combination of the different features, whereas realistic
253 surface current predictions should be based on non linear combinations of features.

254 These non-linear combinations between the different features can instead be incorporated by using deep
255 learning or artificial neural networks. In the following section, we demonstrate the feasibility of using
256 neural networks to extract the nonlinear relationships from data.

5 DEEP LEARNING: ARTIFICIAL NEURAL NETWORKS

257 Artificial neural networks (or neural networks for short) are machine learning algorithms that are loosely
258 modeled after the neuronal structure of a biological brain but on a much smaller scale. A neural network is
259 composed of layers of connected units or nodes called artificial neurons (LeCun et al., 2015; Nielsen, 2015;
260 Goodfellow et al., 2016) that combine input from the data with a set of weights and passes the sum through
261 the node's activation function along with a bias term, to the subsequent set of nodes, to determine to what
262 extent that signal progresses through the network and how it affects the ultimate outcome. Neural nets are
263 typically "feed-forward," meaning that data moves through them in only one direction. A layer is called
264 densely connected when each node in that layer is connected to every node in the layers immediately above
265 and below it. Deep learning, or deep neural networks is the name used for "stacked neural networks" - i.e.,
266 networks composed of several layers.

267 In the past few years, there have been several studies applying machine learning tools, and more
268 specifically deep learning methods to model physical/dynamical processes. For example, deep neural
269 networks (DNN) have been used to develop Reynolds-averaged turbulence models (Ling et al., 2016) to

show that a neural network can be trained to preserve Galilean invariances. Lguensat et al. (2017) developed a convolutional neural network (CNN) based architecture for automated eddy detection and classification from Sea Surface Height (SSH) maps. 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. Pathak et al. (2018) used yet another recently developed machine learning algorithm, known as reservoir computing, to make predictions for the evolution of a very large spatiotemporally chaotic dynamical systems. Another recent study (Bar-Sinai et al., 2018) demonstrated the capabilities of a CNN based method for coarse-graining partial differential equations. This study has strong potential implications for future data-driven subgrid scale parameterizations in atmospheric and oceanographic models. In a recent publication, Gentine et al. (2018) used deep neural networks (DNN), trained with outputs from a superparameterized climate model, to successfully predict most of the key features of embedded convection necessary for climate simulation, thereby suggesting a strong future for data-driven convection parameterizations in climate models. On the oceanographic modeling side, Bolton and Zanna (2018) used CNNs trained on spatio-temporally degraded data from a high-resolution quasi-geostrophic ocean model to successfully replicate the spatio-temporal variability of the eddy momentum forcing. Furthermore, the CNN based method was shown to be generalizable to a range of dynamical behaviours, and could be forced to respect global momentum conservation. One of the common criticisms of deep-learning methods has been that, they are a “black-box”, i.e., lacking any simple intuitive physical interpretations. Some of these recent works (Ling et al., 2016; Bolton and Zanna, 2018; Gentine et al., 2018) showed that data-driven approaches, even with limited data, can be used in conjunction with physical models, to help speed up some of the time intensive/ memory intensive processes in the physical models, while still respecting physical principles.

Our neural network code was written using the Python library Keras (<https://keras.io>) (Chollet et al., 2015), a high-level wrapper around TensorFlow (<http://www.tensorflow.org>). The feed-forward NNs consist of interconnected layers, each of which have a certain number of nodes. The first layer is the input layer, which in our case is a stacked vector containing the input variables just like in the linear regression example above. The last layer is the output layer, which is a stacked vector of the two outputs (U, V). All layers in between are called hidden layers. The activation function, i.e. the function acting on each node – is a weighted sum of the activations in all nodes of the previous layer plus a bias term, passed through a non-linear activation function. For our study, we used the Rectified Linear Unit (ReLU) as an activation function. The output layer is purely linear without an activation function. Training a NN means optimizing the weight matrices and bias vectors to minimize a loss function – in our case the MAE – between the NN predictions and the true values of (u, v) .

The model reduces the loss, by computing the gradient of the loss function with respect to all weights and biases using a backpropagation algorithm, followed by stepping down the gradient – using stochastic gradient descent (SGD). In particular we use a version of SGD called Adam (Kingma and Ba, 2014, 2017). Although most neural network strategies involve normalizing the input variables, we did not use any normalization, since the normalization factors would be largely dependent on the choice of domain / ocean basin, given that the dynamical parameters (like SSH and wind stress) vary widely across the different ocean basins. Instead we wanted the NN to be generalizable across the whole ocean.

We construct a 3-hidden-layer neural network to replace the linear regression model described in the previous section. A schematic model architecture for the neural network is presented in Fig. 2(b). Using the same basic model architecture, we train three NNs on the same three subdomains (Gulf Stream, Kuroshio, ACC) along with one which is trained on the global ocean. Everything including batch size, the training

314 data, the targets, the input features and the number of epochs the model is trained for in each region is kept
315 exactly the same as what we used for the linear regression examples. The only thing that we changed is the
316 model, where instead of 1 layer with no activation we now have three hidden layers with a total of 1812
317 trainable parameters.

318 Just like we did with the linear regression model, we then tracked the evolution of the models' loss
319 function as it moved through batches of input data over multiple epochs (Fig. 3, lighter colored lines in
320 all panels). As we can see, in comparison to the linear regression model, the NNs perform significantly
321 better at reducing the loss in all the ocean regions. What is even more striking is that the NN trained
322 on the globe (dashed line) consistently outperforms the local models, predicting surface currents with
323 lower MAE/ MSE than the models trained on the local subdomains. This is especially noticeable for the
324 Kuroshio region (Fig. 3, second panel), where the NN trained on the globe manages to get the signature of
325 the equatorial currents better than the NN trained specifically in that region (compare panels 3 and 4 of
326 Fig. 5) and gets the absolute error down to $\approx 5\text{cm/s}$. This shows, that in comparison to the linear model
327 the neural network actually manages to learn the physics better when it receives a more spatially diverse
328 input data, and is therefore more generalizable. Even though the linear regression models all manage
329 to get the loss down to comparable magnitude, looking at the spatial plot of the predicted squared error
330 Figs. 45 and 6 gives us an idea how poorly it does at actually learning the physics of surface currents. In
331 comparison, even a relatively shallow 3-hidden-layer neural network performs remarkably better with very
332 few localized hotspots of large errors. This is to be expected since the largest order balance, i.e., geostrophy
333 relies on non-linear combination of the Coriolis parameter (f) with the spatial gradients, and therefore
334 these non-linear combinations are not represented by linear regression and are better captured by a neural
335 network with dense interconnected layers with non-linear activation functions.

336 In Fig. 7 we plot the joint histogram of the zonal and meridional velocity predictions against the true
337 (GCM output) values for the physical model, linear regression model (trained on the local subdomain)
338 and the locally and globally trained neural networks in the ACC sector. From these joint histograms, it
339 is obvious that the physical model, the local and global neural networks all predict velocities that are
340 extremely well correlated with the true velocities in this region. In addition the root mean squared (rms)
341 errors normalized by the rms velocities are also very well correlated between the physical model and neural
342 network predictions. This provides us with reasonable confidence that the model is indeed learning the
343 physics of surface geostrophy and Ekman flow.

344 We also plotted the squared errors in predicted velocity form the physical model (geostrophy+Ekman)
345 and the local Rossby number (expressed as the ratio of the relative vorticity $\zeta = v_x - u_y$, to the planetary
346 vorticity f) in the three domains (Gulf Stream - Fig. 4; Kuroshio - Fig. 5 and the ACC - Fig. 6). It is
347 interesting to note that the localized regions in large root squared errors in both the neural network and
348 physical models coincide with regions where the local Rossby number is high. High Rossby numbers
349 indicate unbalanced flow and the specific regions where we see high Rossby numbers are typically
350 associated with heightened submesoscale activity. We speculate that the prediction errors in these locations
351 are due to the NN's inability to capture higher order balances (e.g. gradient wind, cyclostrophic balance)
352 that are necessary to fully capture the small scale variability associated with these motions and close the
353 momentum budget.

354 The NN also generally predicts weaker velocities near the Equator where the true values of the surface
355 currents are quite large (due to strong Equatorial jets). This can lead to large errors for the global mean,
356 which get magnified when the differences are squared. However we know that geostrophic and Ekman
357 balance also doesn't hold near the Equator. A fairer comparison would therefore involve masking out the

358 near equatorial region ($5^{\circ}N - 5^{\circ}S$) for both the statistical model (i.e. NN predictions) as well as for the
 359 physical model (*geo + ekman*). region sincefails

360 5.1 Neural networks with Convolutional Filters

361 In Section 2 we explained how we can use the local 2×2 stencil to expand the feature vector space by
 362 a factor of 4. We can further expand the feature vector space by passing all the stenciled input features
 363 through k convolutional filters of shape 2×2 . If $k > 4n_f^s$ where n_f^s is the number of input features with a
 364 stencil, we end up with more input features that goes into the NN than before. There is very little functional
 365 difference between this kind of training approach and the one discussed previously, except that we end
 366 up with more trainable parameters, which we can potentially use to extract even more information from
 367 the data. We should note that this is technically not the same as convolutional neural networks, which are
 368 typically used for image analysis and classification, where the convolutional layers serve to *reduce* the
 369 feature vector space without losing information. This is particularly important for problems like image
 370 classification where it is needed to scale down large image datasets without losing feature information. A
 371 schematic of this subcategory of neural network is shown in Fig. 2(c). After applying the convolutional
 372 filter and passing it through a reshape layer in keras the point inputs and filtered stencil inputs are passed
 373 through a Leaky ReLU before being fed into a similar 3-hidden layer NN framework as described before.
 374 Using a similar procedure, we can also apply k 3D convolutional filters of shape $2 \times 2 \times 3$ on the time
 375 and space stenciled inputs to effectively end up with k input features of length N for the stencil variables
 376 (. 2(d)). The goal with the time stenciled input being to potentially learn time derivatives and explore how
 377 the tendencies can affect the NN projections. In hindsight, this data set is probably not be the most suited
 378 for this kind of approach since the variables we used as input features are daily averaged and any fast-time
 379 scale / tendency effects that we hoped to capture from multiple snapshots of the same variable are probably
 380 filtered out by the time averaging. These two approaches are virtually identical with slightly different
 381 preprocessing of the input data.

382 5.2 Model dependence on choice of input features

383 We then trained these NNs with varying combinations of input features to explore how the choice of
 384 input features can influence the model training rate and loss. Feeding the NN models varying combination
 385 of input features, either as stencilled or as point variables and by selectively holding out specific features
 386 for each training case allowed us to assess the relative importance of each physical input variable for the
 387 neural network's predictive capability. The different models with their corresponding input features and the
 388 number of trainable parameters for each case are summarized in table 2. As with all previous examples,
 389 we chose mean absolute error as the loss function for all these experiments. We performed a few training
 390 exercises using the mean squared error instead and did not notice any significant difference. For models
 391 numbered 1 - 13, we used a 2 point space-stencil and for models 1t - 10t, in addition to a stencil in space,
 392 we provide a 3 point time stencil with the intention of helping the neural network "learn" time derivatives.
 393 The different experiments listed in table 2, can broadly be categorized into 6 groups based on their input
 394 features. In group 1, is model 1, where the model only sees η (stencil) and wind stress, τ (point) as input
 395 features. No spatial information is provided. In the second category, we have models that receive η (stencil)
 396 and spatial information X in some form, but no wind stress. This includes models 2, 5t and 7t. The third
 397 category describes models that receive η, θ (stencil) and spatial information X and no wind stress and
 398 includes models 3,6t and 8t. The fourth category describes models that receive SSH (η), spatial information
 399 (X) and wind stress (τ) but no SST and includes models 4, 6, 7, 10, 1t, 3t. The fifth category of models
 400 receive SSH (η), SST (θ), spacial information (X) and wind stress and the only input feature these models
 401 don't receive in any form is sea surface salinity (S). This includes models 5, 8, 9, 11, 2t, 4t. The sixth and

402 final category represents models tat receive all the input features (η, θ, S, X, τ) in some form or another
 403 and includes models 13, 9t and 10t.

404 As mentioned previously, spatial information is provided in one of 3 ways, (a) in the form of 3 dimensional
 405 transformed coordinates (X, Y, Z), (b) just the coriolis parameter (X here serves as a proxy for the coriolis
 406 parameter) and (c) with both the Coriolis parameter and local dx and dy values. Barring a few examples
 407 (models 10, 11) windstress is always provided as a point variable and apart from models 6, 7, 8, 9, none
 408 of the models receive a stencil in the spatial coordinates. We also trained a few models without SSH as
 409 an input feature, but the loss in all these cases was much larger than those shown here ($> 50\text{cm/s}$) and
 410 the NNs fail to pick up any functional dependence on the input features. Those cases are therefore not
 411 presented. Each of these models are trained for 4 Epochs on the same day of data (or 3 consecutive days
 412 centered around that date for the time stencilled cases).

413 In Fig. 9 we summarize the findings from these experiments by plotting the rms error for all the model
 414 predictions along with the rms error for the physical model predictions side by side. With the exception of
 415 5 models (model 1, 5t, 7t, 6t, and 8t) all our NN model predictions have lower domain mean squared errors
 416 than the physics-based models. In terms of features, model without spatial information has the largest error,
 417 followed by models without wind stress (The absolute largest error is for the model without SSH, which is
 418 too big to be considered here). This signifies that to accurately represent surface currents, apart from SSH,
 419 the most important pieces of information required by the neural networks to successfully learn the physics
 420 of surface currents are spatial information and wind stress. It is striking to see how much the model struggles
 421 without spatial information. This implies that latitude dependence is a critical component for a NN to be
 422 able to predict surface currents accurately. It is only expected since the dynamics of surface currents do
 423 depend very strongly on latitude and therefore it is impossible to construct a meaningful prediction model
 424 based on just snapshots without any knowledge of latitude.

425 The zonal mean rms error for the predictions from some of the representative models from the 6 categories
 426 described above are shown in Fig. 8. The NNs all generally predict weaker velocities near the Equator
 427 where the true values of the surface currents are quite large (due to strong Equatorial jets). This can lead
 428 to large errors for the global mean, which get magnified when the differences are squared. However we
 429 know that geostrophic and Ekman balance also doesn't hold near the Equator. Therefore to allow for a
 430 fair comparison between all the models, we mask out the rms errors in a 10° latitude band surrounding the
 431 equator ($5^\circ N - 5^\circ S$) for both the physical and statistical models. Out of all the models, model 1 which
 432 does not receive any spatial information (X), has the highest mean squared errors throughout the globe. For
 433 the models that don't see wind stress (τ) as an input feature, the rms errors are comparable if not smaller at
 434 most latitudes when compared to the physics-based model where you only consider geostrophy (dashed
 435 black line). Additionally, all models that receive η , τ and X in some form perform consistently better than
 436 geostrophy+Ekman at all latitudes (except for near the equator where the physics-based models and the NN
 437 are all equally inadequate). We noticed that during training, the NN's minimize the loss function slightly
 438 faster when a stencil is provided for the spatial coordinates, but after a few epochs the differences in training
 439 loss between models that receive a spatial stencil and models that dont, diminish very rapidly. During
 440 prediction also, the models that receive stencils in spatial coordinates perform slightly better especially at
 441 the high latitudes than the ones where spatial information is provided pointwise.

442 Therefore among the various strategies tested, for this particular dataset, the models that perform the best
 443 in terms of prediction rms error are the models that receive SSH, wind stress and spatial information with a
 444 space stencil. The three point time stencil does not add anything meaningful and appears to hurt, rather
 445 than help the model overall, which was surprising, even though in hindsight we speculate that this might be

446 due to the daily averaged nature of the POP model output. Variables like sea surface temperature and sea
 447 surface salinity have very little impact on the model as well.

448 In terms of choice of features, model 13 stands out as the most practical and physically meaningful
 449 training strategy for a few reasons.

- 450 • It is the most complete in terms of features
 451 • It is the most straightforward to implement, since it does not involve calculating any transformed
 452 3 dimensional coordinates. (All the input variables would be readily available for any gridded
 453 oceanographic dataset.)
 454 • It is one of the models with the lowest prediction rms errors.

455 For these specific reasons we choose model 13 as the reference for performing a sensitivity analysis. The
 456 purpose of this analysis is to characterize the sensitivity of the model to perturbations in the different
 457 input features during testing/prediction. For the sensitivity tests we simply add a gaussian noise of varying
 458 amplitude to each of the input variables, while keeping the rest of the input variables fixed. For each
 459 of the input variables ($x_i \in \{\eta, \theta, X, dx, dy, \tau_x, tau_y\}$), we chose 3 different zero-mean gaussian noise
 460 perturbations with the standard deviations of $0.5\sigma(x_i)$, $\sigma(x_i)$, and $2\sigma(x_i)$, where $\sigma(x_i)$ is the standard
 461 deviation of the corresponding input variable x_i . The model loss is then evaluated for each of these
 462 perturbations and normalized by the amplitude of perturbations (right panel Fig. 10). This normalization is
 463 done to level the playing field for all the input variables and allow for a one-to-one comparison since the
 464 different input variables vary in orders of magnitude (e.g. the amplitude of perturbations in SSH is $O(100)$,
 465 while the amplitude of perturbations in wind stress is $O(1)$ and therefore a perturbation of amplitude $\sigma(\eta)$
 466 in η would lead to a much larger model error than a perturbation of $\sigma(\tau_x)$ in τ_x would, as can be seen from
 467 the log scaling of the y-axis in the left panel of Fig. 10).

468 Given what we learned about the importance of spatial coordinates for NN training, it is not surprising
 469 to see that for prediction also, the NN is most sensitive to perturbations in the coriolis parameter (or X).
 470 The input variables that the model is most sensitive to, arranged in descending order of model sensitivity
 471 are coriolis parameter, SSH and wind stress, followed by SST. The model is not particularly sensitive to
 472 perturbations in local grid spacing or salinity. The relative effect of the input variables, observed in the
 473 model sensitivity test closely matches what we saw in the different model training examples where we
 474 selectively held out these features. This again confirms that in order to train a deep learning model to make
 475 physically meaningful and generalisable predictions of surface currents it is not sufficient to simply provide
 476 it snapshots of dynamical variables like SSH as images. We also need to provide spatial information like
 477 latitude for the NN's to effectively "learn" the physics of surface currents.

6 SUMMARY AND FUTURE DIRECTIONS

478 The goal of this study was to use machine learning to make predictions of ocean surface currents from
 479 satellite observable quantities like SSH wind stress, SST etc. Our central question was: Can we train deep
 480 learning based models to learn physical models of estimating surface currents like geostrophy, Ekman flow
 481 and perhaps do better than the physical models themselves?

482 We used the output from the CESM POP model surface fields as our "truth" data for this study. As a
 483 first order example, we tested a linear regression model for a few of local subdomains extracted from
 484 the global GCM output. Linear regression works well only when the domains are small and far removed
 485 from the equator and gets progressively worse as the domain gets bigger and the variation in local coriolis
 486 parameter gets large. It performs most poorly when f changes sign in the domain. reasonably well for

487 small enough regions, far enough from the equator. This showed that unsurprisingly it is not possible to
488 train a simple linear model to accurately predict surface currents. In addition, providing more data does not
489 necessarily improve the predictive ability of a linear model and only made it worse as it starts overfitting.
490 Whereas for the same kind of domain, a neural network we can minimize the loss (MSE) with fewer
491 data-points and still remain generalisable, since neural networks can learn functional relationships between
492 regressors (input features) with only a small amount of data. The model's ability to make predictions is also
493 shown to improve with more data. Furthermore, compared to a linear regression model, a NN even with a
494 relatively small network of densely connected nodes, with a suitable non-linear activation function (like
495 ReLU), allows us to have a large number of trainable parameters (weights, biases) that can be optimized to
496 minimize the loss. The activation function is what allows the different non-linear combinations between the
497 different regressors (input features). Furthermore, a neural network trained on the entire globe is shown to
498 predict surface currents more accurately in the sub-domains than neural networks trained in those specific
499 sub domains. In comparison, a similar approach with a linear regression model produces the opposite
500 result, *i.e.* a globally trained linear regression model produces higher prediction errors than the one that's
501 trained on each specific sub-domain. The fact that spatially diverse data actually makes the neural network
502 perform better is indicative of the fact that a neural network can actually "learn" the functional relationships
503 needed for calculating surface currents rather than simply memorizing some target values for different
504 combination of input features. By examining the dependence of the NNs on the choice of input features
505 and by looking at the sensitivity of a NN model to perturbations in the input features, we established that
506 apart from SSH, the physical location of the input features is one of the most crucial elements for the
507 NN to "learn" the physics of surface currents. It is further demonstrated that with a careful and deliberate
508 choice of input features the neural network can even beat the physics-based models at predicting surface
509 currents accurately in most regions of the global ocean. A key ingredient for calculating the Ekman part of
510 the flow using current physics based models is the vertical diffusivity, which is largely unknown for most
511 of the global ocean. Most observational ocean current estimates that include the Ekman part of the flow
512 relies on inferring the vertical diffusivity based on empirical multiple linear regressions with Lagrangian
513 surface drifter data. The neural network approach, by comparison does not suffer from the same kind of
514 limitation, since in this framework, we do not need to provide A_z as an input feature, which is one more
515 added advantage for this method.

516 In this study, we wanted to see whether we can train a statistical model like a NN with data to essentially
517 match or perhaps beat the baseline physics based models we currently use to estimate surface currents. By
518 examining the errors in surface current predictions from our NN predictions and comparing them with
519 predictions from physically motivated models (like geostrophy and Ekman dynamics), we showed that
520 a relatively simple NN captures most of the large scale flow features equally well if not better than the
521 physical models, with only one day of training data for the globe.

522 However, some key aspects of the flow, associated with mesoscale and sub-mesoscale turbulence are not
523 reproduced. We speculate that this is possibly caused by the fact that the neural network framework can
524 not capture the higher order balances (gradient wind) that are likely at play in these regions since these
525 hotspots of high errors are collocated with regions of High Ro where balance breaks down (see Figs. 4-6).

526 One of the biggest hurdles associated with these studies is figuring out efficient strategies to stream large
527 volumes of earth system model data into a NN framework. So before diving headfirst into the highest
528 resolution global ocean model (currently available), we wanted to test the feasibility of using a regression
529 model based on deep learning as a framework for estimating surface currents with a lower resolution
530 model data (smaller/more manageable dataset), while still being eddy resolving. Hence we chose the CESM

531 POP model data for this present study. In the future, we propose to train a NN with data from a higher
532 spatio-temporal resolution global ocean model like the MITgcm llc4320 model (Rocha et al., 2016). As a
533 further step, we could coarse-grain such a model to SWOT-like resolutions, or use the SWOT simulator,
534 train NNs on that, and make predictions for global surface currents.

535 As for the weak surface currents predicted by our NN at the equator, we need to keep in mind that
536 geostrophic balance (defined by the first order derivatives of SSH) only holds away from the equator
537 and satellite altimetry datasets (e.g. AVISO, Ducet et al., 2000) typically employ a higher order balance
538 (Lagerloef et al., 1999) at the equator, to match the flow regime with the geostrophic regime away from the
539 equator. One possible way to train the NN to learn these higher order balances would be by increasing the
540 stencil size around each point. Since our primary goal was for the model to learn geostrophy, we started
541 with a spatial stencil in SSH. We also explored training approaches where we provided stencils in SST, wind
542 and SSS, with the intention of helping the model learn about wind-stress curl and thermal wind balance. In
543 practice, however these didn't payoff as much and these additional stencils did not significantly improve
544 model performance. In future approaches we can try to provide separate stencils of varying size to each of
545 these input variables, to test whether we can further improve the model accuracy.

546 As another future step, we also aim to incorporate recursive neural networks (RNNs) in conjunction
547 with convolutional filters of varying kernel sizes, to train the models on cyclostrophic or gradient wind
548 balance. This recursive neural network approach would be analogous to iteratively solving the gradient
549 wind equation (Knox and Ohmann, 2006), a technique which was originally developed for numerical
550 weather prediction before advances in computing allowed for integrating the full non-linear equations.

551 The present work demonstrates that to a large extent, a simple neural network can be trained to extract
552 functional relationships between SSH, wind stress etc. and surface currents with quite limited data. The
553 field of deep learning as of now is rapidly evolving. It remains to be seen, if with some clever choices of
554 training strategies and by using some of the other more recently developed deep learning techniques, we
555 can improve upon this. In this study, we propose a few approaches that can be implemented to improve
556 upon our current results and would like to investigate this in further detail in future studies. In addition,
557 we believe that data driven approaches, like the one shown in this present study, have strong potential
558 applications for various practical problems in physical oceanography, and require further exploration.
559 Insights gained from this type of analysis could be of great potential significance, especially for future
560 satellite altimetry missions like SWOT.

CONFLICT OF INTEREST STATEMENT

561 The authors declare that the research was conducted in the absence of any commercial or financial
562 relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

563 The authors confirm contribution to the paper as follows: study conception and design: AS, RPA; training
564 and testing of statistical models: AS ; analysis and interpretation of results: AS ; draft manuscript
565 preparation: AS, RPA. All authors reviewed the results and approved the final version of the manuscript.

FUNDING

566 The authors acknowledge support from NSF Award OCE 1553594 and NASA award NNX16AJ35G
567 (SWOT Science Team). The work was started in 2018 and an early proof-of-concept was reported in AS's
568 PhD dissertation (Sinha, 2019) .

ACKNOWLEDGMENTS

569 The authors acknowledge support from NSF Award OCE 1553594 and NASA award NNX16AJ35G
570 (SWOT Science Team).

DATA AVAILABILITY STATEMENT

571 The datasets used for this study can be found in the PANGEA data catalog on Google cloud storage under
572 and all the relevant code used for all the analysis presented here can be found in the .

REFERENCES

- 573 Arbic, B. K., Richman, J. G., Shriver, J. F., Timko, P. G., Metzger, E. J., and Wallcraft, A. J. (2012).
574 Global modeling of internal tides: Within an eddying ocean general circulation model. *Oceanography*
575 25, 20–29
- 576 Bar-Sinai, Y., Hoyer, S., Hickey, J., and Brenner, M. P. (2018). Data-driven discretization: a method for
577 systematic coarse graining of partial differential equations. *ArXiv e-prints*
- 578 [Dataset] Bolton, T. and Zanna, L. (2018). Applications of deep learning to ocean data inference and
579 sub-grid parameterisation. doi:10.31223/osf.io/t8uhk
- 580 Bonjean, F. and Lagerloef, G. S. E. (2002). Diagnostic model and analysis of the surface currents in the
581 tropical pacific ocean. *Journal of Physical Oceanography* 32, 2938–2954. doi:10.1175/1520-0485(2002)
582 032<2938:DMAAOT>2.0.CO;2
- 583 Chapman, C. and Charantonis, A. A. (2017). Reconstruction of subsurface velocities from satellite
584 observations using iterative self-organizing maps. *IEEE Geoscience and Remote Sensing Letters* 14,
585 617–620. doi:10.1109/LGRS.2017.2665603
- 586 [Dataset] Chollet, F. et al. (2015). Keras. <https://keras.io>
- 587 Ducet, N., Le Traon, P. Y., and Reverdin, G. (2000). Global high-resolution mapping of ocean circulation
588 from topex/poseidon and ers-1 and -2. *Journal of Geophysical Research: Oceans* 105, 19477–19498.
589 doi:10.1029/2000JC900063
- 590 Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., and Yacalis, G. (2018). Could machine learning break
591 the convection parameterization deadlock? *Geophysical Research Letters* 45, 5742–5751. doi:10.1029/
592 2018GL078202
- 593 Goodfellow, I., Bengio, Y., Courville, A., and Bengio, Y. (2016). *Deep learning*, vol. 1 (MIT press
594 Cambridge)
- 595 Gregor, L., Kok, S., and Monteiro, P. M. S. (2017). Empirical methods for the estimation of southern
596 ocean CO_2 : support vector and random forest regression. *Biogeosciences* 14, 5551–5569. doi:10.5194/
597 bg-14-5551-2017
- 598 Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *CoRR* abs/1412.6980
- 599 [Dataset] Kingma, D. P. and Ba, J. (2017). Adam: A method for stochastic optimization
- 600 Knox, J. A. and Ohmann, P. R. (2006). Iterative solutions of the gradient wind equation. *Computers
601 Geosciences* 32, 656–662. doi:<https://doi.org/10.1016/j.cageo.2005.09.009>
- 602 Lagerloef, G. S. E., Mitchum, G. T., Lukas, R. B., and Niiler, P. P. (1999). Tropical pacific near-surface
603 currents estimated from altimeter, wind, and drifter data. *Journal of Geophysical Research: Oceans* 104,
604 23313–23326. doi:10.1029/1999JC900197
- 605 LeCun, Y., Bengio, Y., and Hinton, G. (2015). Deep learning. *nature* 521, 436
- 606 LeTraon, P. and Morrow, R. (2001). Ocean currents and eddies 69, 171–215
- 607 Lguensat, R., Sun, M., Fablet, R., Mason, E., Tandeo, P., and Chen, G. (2017). Eddynet: A deep neural
608 network for pixel-wise classification of oceanic eddies. *CoRR* abs/1711.03954

- 609 Ling, J., Kurzawski, A., and Templeton, J. (2016). Reynolds averaged turbulence modelling using deep
610 neural networks with embedded invariance. *Journal of Fluid Mechanics* 807, 155–166. doi:10.1017/jfm.
611 2016.615
- 612 Morrow, R., Blurmstein, D., and Dibarboore, G. (2018). Fine-scale altimetry and the future swot mission.
613 *New Frontiers in Operational Oceanography*, 191–226
- 614 Nielsen, M. A. (2015). *Neural networks and deep learning*, vol. 25 (Determination press USA)
- 615 Pathak, J., Hunt, B., Girvan, M., Lu, Z., and Ott, E. (2018). Model-free prediction of large spatiotemporally
616 chaotic systems from data: A reservoir computing approach. *Phys. Rev. Lett.* 120, 024102. doi:10.1103/
617 PhysRevLett.120.024102
- 618 Rocha, C. B., Chereskin, T. K., Gille, S. T., and Menemenlis, D. (2016). Mesoscale to submesoscale
619 wavenumber spectra in drake passage. *Journal of Physical Oceanography* 46, 601–620
- 620 Sinha, A. (2019). *Temporal Variability in Ocean Mesoscale and Submesoscale Turbulence*. Ph.D. thesis,
621 Columbia University in the City of New York. doi:<https://doi.org/10.7916/d8-bngk-r215>
- 622 Small, R. J., Bacmeister, J., Bailey, D., Baker, A., Bishop, S., Bryan, F., et al. (2014). A new synoptic scale
623 resolving global climate simulation using the community earth system model. *Journal of Advances in
624 Modeling Earth Systems* 6, 1065–1094. doi:<https://doi.org/10.1002/2014MS000363>
- 625 Smith, R., Jones, P., Briegleb, B., Bryan, F., Danabasoglu, G., Dennis, J., et al. (2010). The parallel ocean
626 program (pop) reference manual ocean component of the community climate system model (ccsm) and
627 community earth system model (cesm). *Rep. LAUR-01853* 141, 1–140
- 628 Sudre, J., Maes, C., and Garcon, V. (2013). On the global estimates of geostrophic and ekman
629 surface currents. *Limnology and Oceanography: Fluids and Environments* 3, 1–20. doi:10.1215/
630 21573689-2071927
- 631 Sudre, J. and Morrow, R. A. (2008). Global surface currents: a high-resolution product for investigating
632 ocean dynamics. *Ocean Dynamics* 58, 101. doi:10.1007/s10236-008-0134-9
- 633 Uchida, T., Abernathey, R., and Smith, S. (2017). Seasonality of eddy kinetic energy in an eddy permitting
634 global climate model. *Ocean Modelling* 118, 41 – 58. doi:<https://doi.org/10.1016/j.ocemod.2017.08.006>

Table 1. Table summarizing model errors from the the physics based model (geostrophy + Ekman flow) and the two types of regression models - linear regression and neural network (Panel (a) and (b) in Fig. 2).

Model (training region)	Number of trainable parameters	Epochs	MAE (train) [cm/s]	MAE (eval) GS [cm/s]	MAE (eval) Kuroshio [cm/s]	MAE (eval) ACC [cm/s]
LR (Gulf Stream)	38	8	10.7	11.4	-	-
NN (GS)	1812	8	2.3	3.7	-	-
LR (Kuroshio)	38	5	12.9	-	13.4	-
NN (Kuroshio)	1812	5	5.8	-	7.0	-
LR (ACC)	38	5	7.5	-	-	7.5
NN (ACC)	1812	5	1.9	-	-	4.5
NN (global)	1812	4	3.0	2.4	5.1	1.8
<i>geo+Ek</i> (global)	-	-	-	6.1	29.2	3.9

Table 2. Table summarizing the different CNNs and the training strategies explored

Model No.	Stencil in space (2s)	Stencil in time (3t)	Stencil Variables	Point Variables	Number of trainable parameters
1	✓	✗	η	τ_x, τ_y	4772
2	✓	✗	η	$X (= \frac{f}{2\Omega})$	4732
3	✓	✗	η, θ	X	5052
4	✓	✗	η	X, τ_x, τ_y	4812
5	✓	✗	η, θ	X, τ_x, τ_y	5132
6	✓	✗	η, X, Y, Z	τ_x, τ_y	5732
7	✓	✗	η, X	τ_x, τ_y	5092
8	✓	✗	η, θ, X, Y, Z	τ_x, τ_y	6052
9	✓	✗	η, θ, X	τ_x, τ_y	5412
10	✓	✗	η, τ_x, τ_y	X	5372
11	✓	✗	$\eta, \theta, \tau_x, \tau_y$	X	5692
12	✓	✗	η, θ	$X, dx, dy, \tau_x, \tau_y$	5212
13	✓	✗	η, θ, S	$X, dx, dy, \tau_x, \tau_y$	5532
1t	✓	✓	η	$\tau_x, \tau_y, X, dx, dy$	5532
2t	✓	✓	η, θ	$\tau_x, \tau_y, X, dx, dy$	6492
3t	✓	✓	η	τ_x, τ_y, X	5452
4t	✓	✓	η, θ	τ_x, τ_y, X	6412
5t	✓	✓	η	X, dx, dy	5452
6t	✓	✓	η, θ	X, dx, dy	6412
7t	✓	✓	η	X	5372
8t	✓	✓	η, θ	X	6332
9t	✓	✓	η, θ, S	τ_x, τ_y, X	7372
10t	✓	✓	η, θ, S	$\tau_x, \tau_y, X, dx, dy$	7452

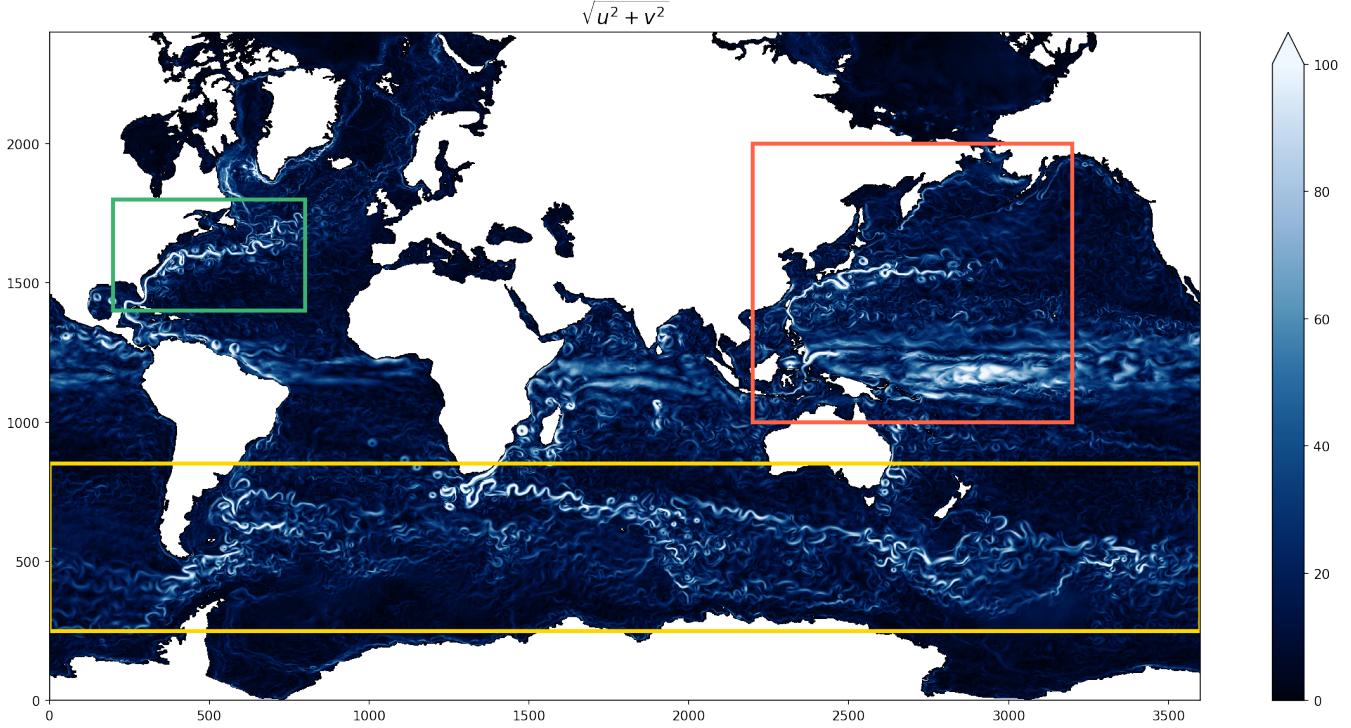


Figure 1. Snapshot of the surface speed in the CESM POP model with the 3 boxes in different colors indicating the training regions chosen for the different regression models. The green box is chosen as the Gulf stream region, the red box is Kuroshio and the yellow box represents the Southern Ocean / Antarctic circum polar current (ACC). The Kuroshio region extends slightly south of the equator to include the equatorial jets in the domain and to test the models' ability to generalize to large variations in f .

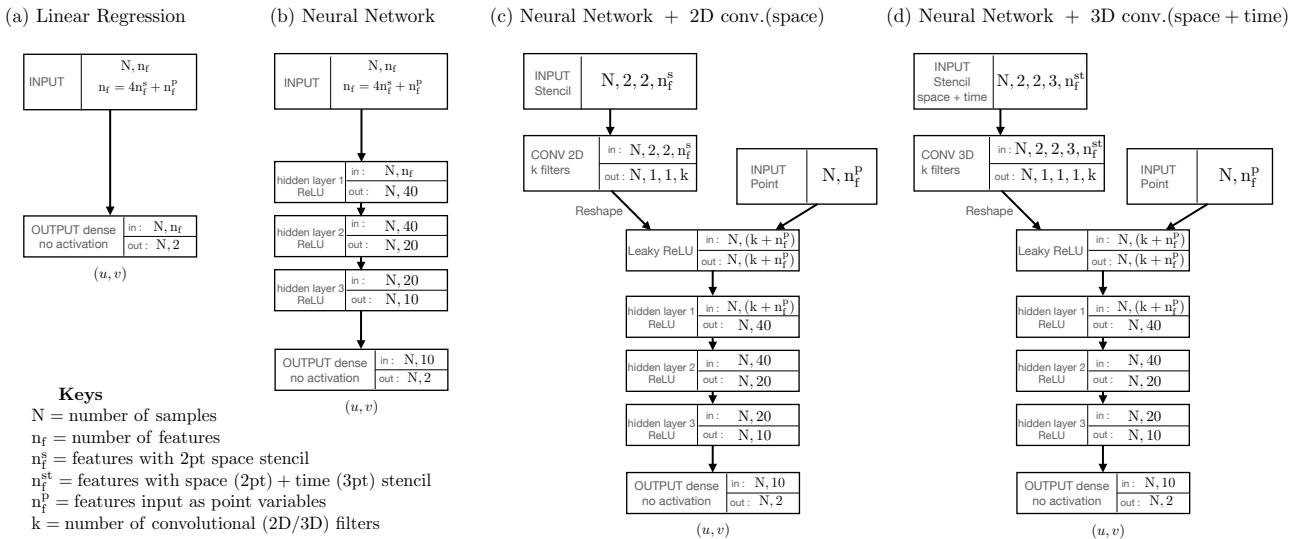


Figure 2. Schematic of the 4 different types of statistical models used in the study. All models shown were implemented using keras tensorflow (Chollet et al., 2015) and we use Mean absolute error (MAE) as the loss function and the Adam optimizer Kingma and Ba (2017) with default parameters and learning rates.

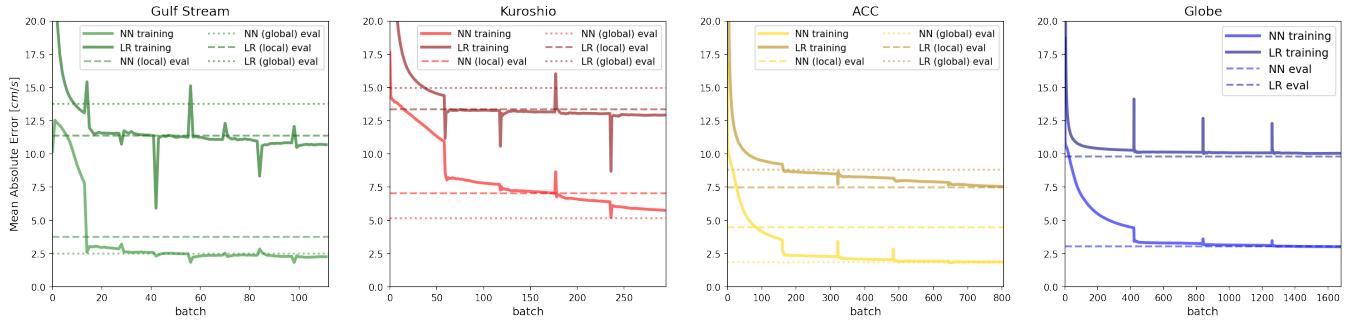


Figure 3. Evolution of the loss function (mean absolute error; MAE) for Neural Networks and Linear regression models during training. Horizontal lines of the corresponding color denote the MAE for the model when evaluated at a different time snapshot. Dashed lines denote the evaluated (test data) MAE for the local model and dotted lines denote that for the model trained on the globe.

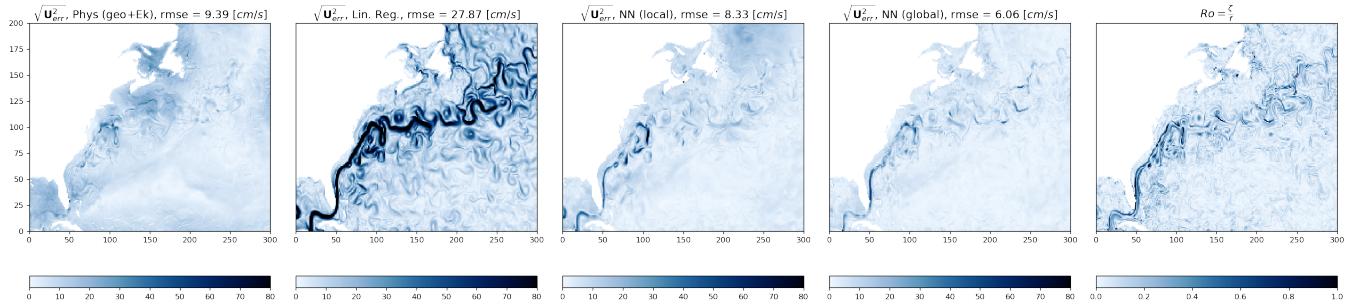


Figure 4. Snapshot of model predicted root square errors for the physics based model (left) and the 3 different regression models - Linear regression (second from left), neural network, trained on this local domain (third panel) and neural network, trained on the globe (4th panel) compared side by side with the local Rossby Number (Ro, right panel) in the Gulf Stream region indicated by the green box in Fig. 1.

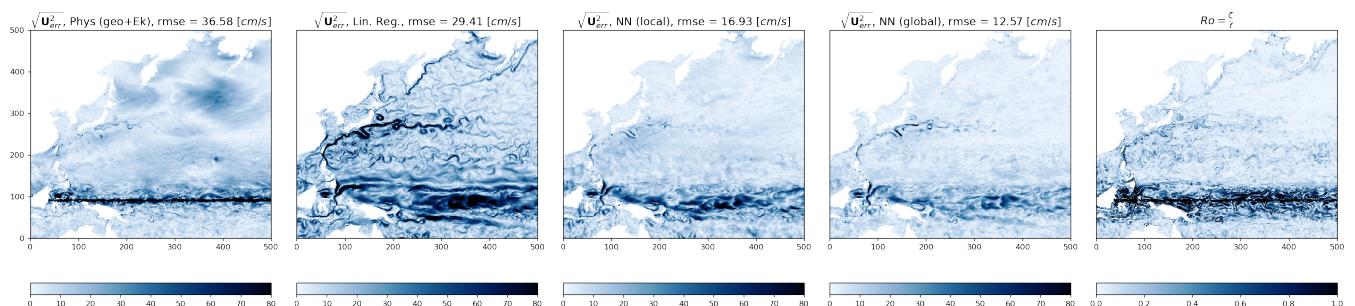


Figure 5. Snapshot of model predicted root square errors for the physics based model (left) and the 3 different regression models - Linear regression (second from left), neural network, trained on this local domain (third panel) and neural network, trained on the globe (4th panel) compared side by side with the local Rossby Number (Ro, right panel) in the Kuroshio region indicated by the red box in Fig. 1. Note the large errors in all the model predictions near the equator.

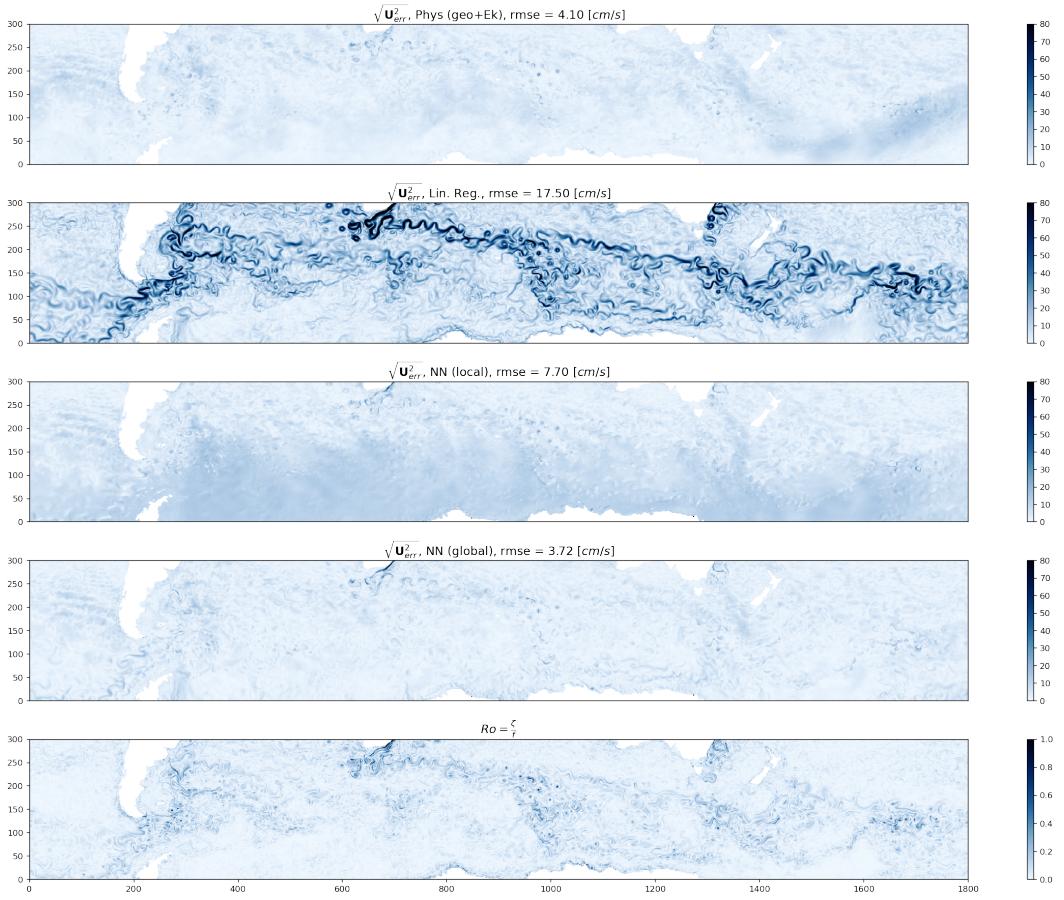


Figure 6. Snapshot of model predicted root square errors for the physics based model (top) and the 3 different regression models - Linear regression (second panel), neural network, trained on this local domain (third panel) and neural network, trained on the globe (4th panel) compared side by side with the local Rossby Number (Ro, bottom panel) in the Southern Ocean/ Antarctic circumpolar current region indicated by the yellow box in Fig. 1.

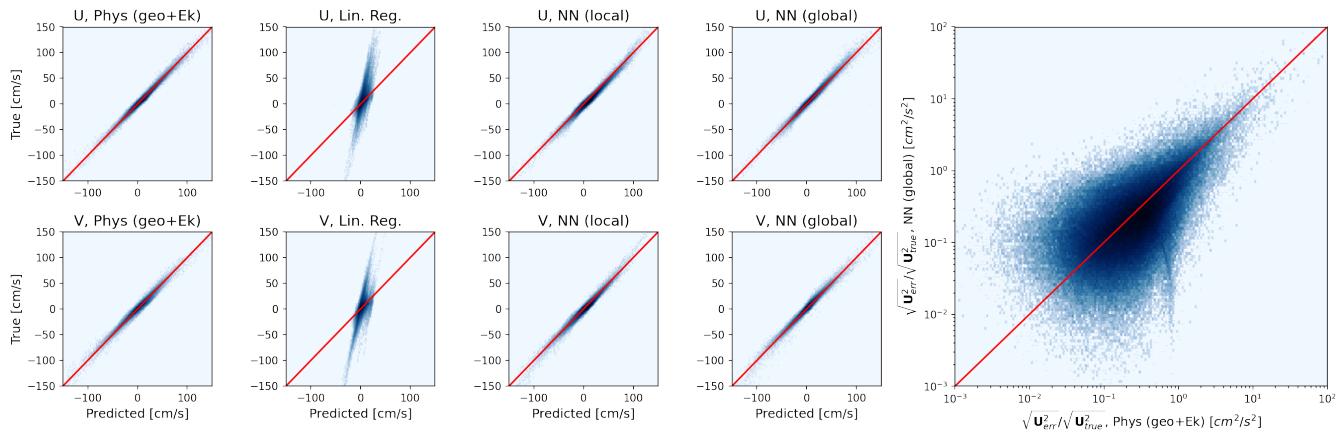


Figure 7. Scatterplot of true v predicted zonal and meridional velocities for the different physical and regression models (8 panels on the left) in the ACC region. The right panel shows the scatterplot of the root mean squared errors (normalized by the root mean square velocities) for the physical and neural network model predictions.

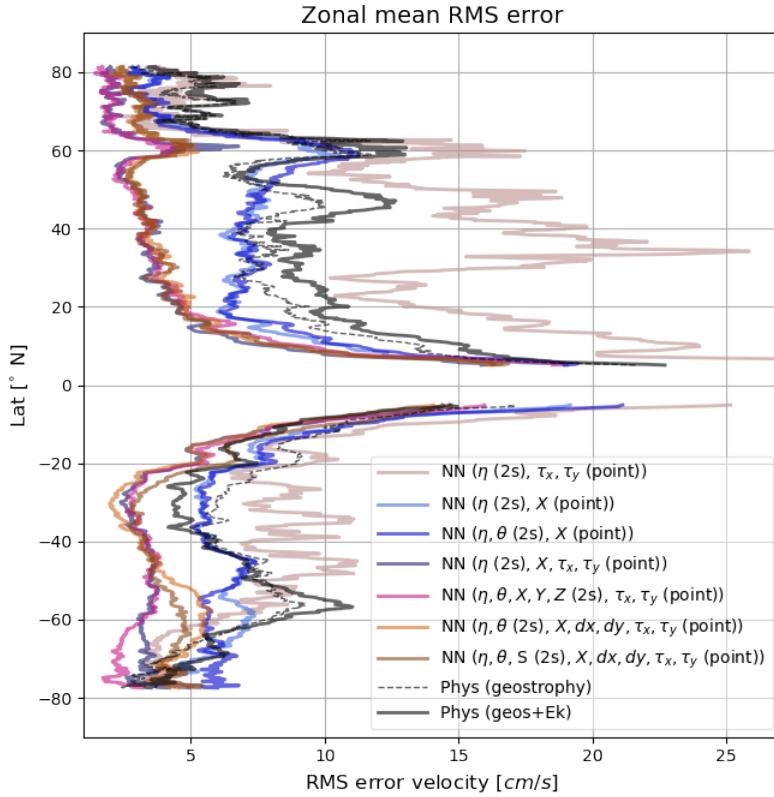


Figure 8. Comparison of the zonal mean rms errors for the various NN predictions shown alongside the physical model (with and without Ekman flow).

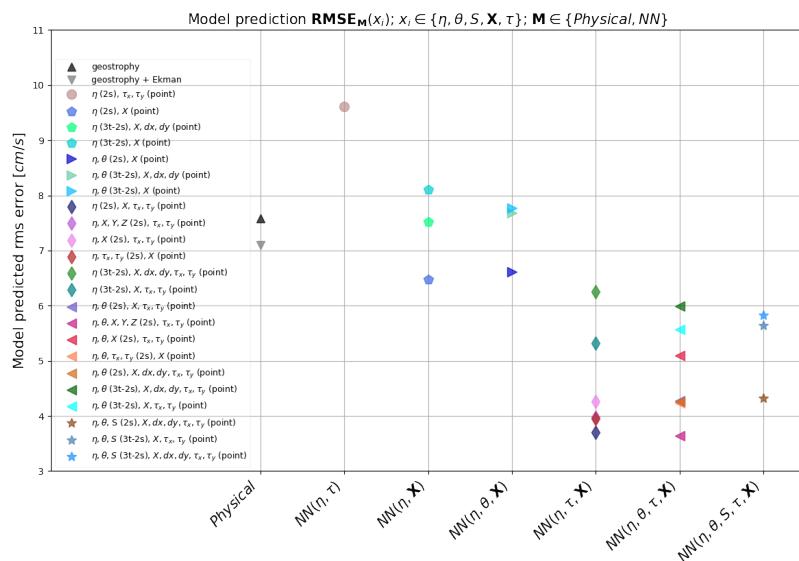


Figure 9. Figure comparing the rms error of the different model predictions along with the rms error for the physical models as a function of features.

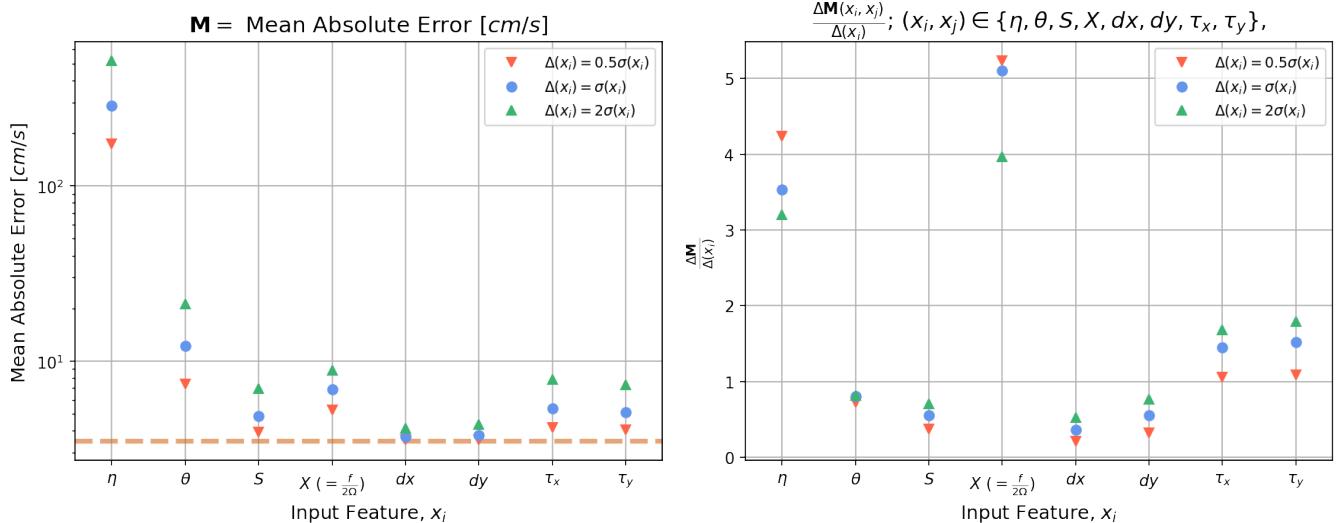


Figure 10. Sensitivity of the neural networks to perturbations in the different input features. Each input feature is perturbed by 3 different gaussian noise perturbations with standard deviations of 0.5σ , σ , and 2σ , where σ is the standard deviation of each variable, while keeping the remaining input variables fixed. The left panel shows the model loss (mean absolute error, MAE) evaluated for each of these perturbations. The horizontal dashed line represents the loss for the unperturbed/control case. The right panel shows the deviation in MAE for each of these perturbation experiments normalized by the amplitude of the perturbation.