

1 **Design and implementation of a data-driven**
2 **parameterization for mesoscale thickness fluxes**

3 **Dhruv Balwada¹, Pavel Perezhogin², Alistair Adcroft⁴, and Laure Zanna^{2,3}**

4 ¹Lamont-Doherty Earth Observatory, Columbia University, Palisades, NY, USA

5 ²Courant Institute of Mathematical Science, New York University, New York, NY, USA

6 ³Center for Data Science, New York University, New York, NY, USA

7 ⁴Princeton University, Princeton, NJ, USA

8 **Key Points:**

- 9 • A data-driven mesoscale thickness flux parameterization was designed and imple-
10 mented in MOM6 as an alternative to the Gent-McWilliams scheme
11 • The parameterization is stable and skillful across a range of resolutions including
12 the eddy-permitting gray zone
13 • The parameterization is able to reduce potential energy without overly dissipat-
14 ing the resolved eddy energy

Abstract

Mesoscale eddies are a major sink of available potential energy (APE) in the ocean. When these eddies are not resolved or only partially resolved in a model, this effect needs to be parameterized to simulate a realistic large-scale state. Traditionally, the Gent-McWilliams (GM) parameterization has provided this sink of APE. However, the GM parameterization, which diffuses isopycnal heights, is not accompanied by a skillful prescription for GM diffusivity rooted in data from observations or models. Also, at eddy-permitting resolutions, GM diffusion can negatively impact resolved eddies, and the only scale-aware prescription is to turn GM off in regions where eddies are permitted. Here we present a novel data-driven parameterization, as a substitute for GM, that extracts APE without overly negative impacts on the resolved flow. It is both flow-aware and scale-aware, and its magnitude can be tuned using an $O(1)$ non-dimensional number. Features like non-dimensional inputs/outputs, lateral non-locality, flow-dependent coordinates, and range limitations improve the generalization of the data-driven scheme. Functional forms are learned via a small multi-layer perceptron, ensuring low computational cost and simple implementation in ocean models. The parameterization performs skillfully in offline evaluation, especially at scales smaller than the largest eddies. It is implemented in NOAA GFDL's MOM6 and shown to be skillful in online tests in two-layer idealized simulations of a zonal channel and wind-driven gyre, at both eddy-permitting and non-eddy-permitting resolutions. This work suggests a path towards leveraging high-resolution simulations for the reduction of structural error and improvement in the fidelity of climate simulations.

Plain Language Summary

Mesoscale (100 km) eddies are the dominant flows in the ocean and play a key role in shaping large-scale circulation features such as wind-driven gyres and the meridional overturning circulation. Since these eddies are not fully resolved in many modern ocean models – especially those that are run for long periods or include many ensemble members – their effects must be represented through parameterizations. A commonly used approach, the Gent-McWilliams (GM) parameterization, removes available potential energy (APE) from the system but lacks a data-driven way to set its strength. Moreover, at eddy-permitting resolutions, GM can interfere with the resolved flows.

We present a new data-driven parameterization designed to better represent eddy effects in ocean models. It learns a functional form from high-resolution simulations using a compact neural network and is designed to be flow-aware, scale-aware, and computationally efficient. The parameterization is implemented in the MOM6 ocean model and shows skillful performance both offline and in idealized simulations, especially at scales smaller than the largest eddies. It extracts APE from large scale flow without degrading resolved features, offering a promising alternative to GM for a wide range of ocean modeling applications.

1 Introduction

Ocean circulation models solve equations describing the motions in the ocean on a finite-size discrete grid, and are thus unable to resolve the phenomena at scales smaller than the grid scales. However, it is often the case that these *sub-grid* phenomena are more than just small-scale variability, and could have a profound impact on the characteristics of the resolved motions. To ensure fidelity of the model behavior at the resolved scales, the effects of these sub-grid phenomena need to be appropriately *parameterized*. One important and often unresolved range of scales in ocean models are the mesoscales.

Mesoscales ($\sim 50\text{--}200\text{km}$) are the dominant energy-containing scales in the ocean (Ferrari & Wunsch, 2009), and subsequently play an important role in shaping the mean

64 circulation and stratification (Gent, 2011), and in transporting tracers (Abernathey &
65 Wortham, 2015). These eddies are believed to be largely generated as a result of baro-
66 clinic instability (K. S. Smith, 2007), which has its fastest growth rate at scales close to
67 the first baroclinic deformation radius (scales of $\sim 10 - 100\text{km}$) (Chelton et al., 1998;
68 Tulloch et al., 2011). These instabilities have a tendency to grow at the expense of the
69 available potential energy (APE), and thus their bulk effect is to flatten isopycnals in
70 the ocean. The dominance of rotation at these scales also usually leads to a subsequent
71 inverse energy cascade, which is why the dominant peak of energy is usually larger than
72 the dominant scale of the instability (Tulloch et al., 2011). Thus, to properly resolve these
73 eddies and their effects the model grid spacing needs to be at least as small as the de-
74 formation radius (Hallberg, 2013). The inverse cascade also takes place in the vertical
75 (K. S. Smith & Vallis, 2002), transferring energy to the barotropic and first couple of baro-
76 clinic modes, and the associated flows tend to have weak vertical shear and do not gen-
77 erate of small-scale turbulence or diapycnal mixing — mesoscale eddies are dominantly
78 adiabatic processes in the interior of the ocean.

79 Conventionally, mesoscale eddy effects have been parameterized using the Gent-
80 McWilliams (GM) parameterization (Gent & McWilliams, 1990), particularly in mod-
81 els that do not resolve the deformation radius. This parameterization was designed to
82 respect two important aspects of mesoscale processes: (i) the parameterization is adi-
83 abatic, and (ii) the net effect of the parameterization is to reduce available potential en-
84 ergy. In models with depth as the vertical coordinate, this is achieved by representing
85 the horizontal eddy fluxes in the form of a horizontally downgradient buoyancy diffu-
86 sion and then setting the vertical component of the eddy flux to be upgradient, such that
87 the net flux is along isopycnals and the resulting operator behaves like advection. In con-
88 trast in isopycnal models, this is achieved by diffusing the interface heights, which can
89 also be represented as extra eddy driven advection. These recipes led to dramatic qual-
90 itative improvements in the ocean simulations, particularly the adiabatic aspect ensures
91 that water masses were not eroded away in the interior by spurious diffusion. However,
92 the associated eddy diffusivity has always been a major source of uncertainty and a very
93 active topic of research for decades (e.g. Visbeck et al., 1997; Ferreira et al., 2005; Eden
94 & Greatbatch, 2008; Marshall et al., 2012; Jansen et al., 2015). Also, as with all con-
95 ventional ocean parameterizations, this scheme was designed to represent the bulk ef-
96 fects of eddies and not to have any skill in representing the spatial or temporal struc-
97 tures of the eddy effects, which leads to detrimental effects at resolutions where eddies
98 are partially resolved (Hallberg, 2013; Mak et al., 2023).

99 In recent years, machine learning (ML) based methods have started to show a lot
100 of promise in improving different aspects of computational modeling, including improv-
101 ing parameterizations (Bracco et al., 2025; Lai et al., 2024). While conventional param-
102 eterizations require domain scientists to develop mathematical operators that are at best
103 able to mimic the bulk effects of sub-grid phenomena, ML methods directly learn the
104 functional relationship between the sub-grid effects and the large-scale fields using ap-
105 propriate data. In many instances it is also possible to design the ML models to obey
106 some physical properties. These methods have led to the development of parameteriza-
107 tions for idealized systems (Ross et al., 2023; Srinivasan et al., 2023), thermodynamic
108 and momentum tendencies in the atmosphere (Brenowitz & Bretherton, 2018; Yuval et
109 al., 2021), boundary layers processes in the ocean (Sane et al., 2023; Ramadhan et al.,
110 2020; Bodner et al., 2023), and momentum tendencies in the ocean (Zhang et al., 2023;
111 Zanna & Bolton, 2020; Perezhogin et al., 2024). All these new parameterizations have
112 shown improved skill over conventional parameterizations, and the potential issues raised
113 about implementation, stability and generalization are rapidly being addressed. While
114 some sub-grid effects of ocean eddies have been investigated using this approach, the im-
115 pact of mesoscale eddies on the density or thickness field - the aspect parameterized by
116 the GM parameterization - has not yet been addressed using data-driven parameteriza-
117 tions.

118 Here we present a new data-driven parameterization to account for the impact that
 119 mesoscale eddies have on the thickness field in the ocean. Our parameterization is de-
 120 signed to ensure that it maintains the important adiabatic constraint that was introduced
 121 by the GM parameterization. However, unlike GM, our parameterization has not been
 122 designed to be a local sink of APE; rather it is optimized to capture the spatial struc-
 123 ture of the eddy effects. With an eye towards implementation, we use a small fully con-
 124 nected neural network - multi-layer perceptron (MLP), which can be easily added to any
 125 ocean model code. Here we implemented this into GFDL’s Modular Ocean Model 6 (MOM6).
 126 In section 2 we describe the filtering framework that is used to diagnose the sub-grid im-
 127 pact of eddies from high resolution (HR) simulations, and in section 3 describe how to
 128 cast these sub-grid effects into a functional form that potentially has some ability to gen-
 129 eralize to unseen data. Also, in section 3 we describe the machine learning architecture,
 130 and training process. In section 4 we describe the HR datasets and how they were pro-
 131 cessed. In section 5 we show that our parameterization is successful in both an offline
 132 and online sense, and finally, in section 6 we conclude with a discussion, potential caveats
 133 and an outlook towards future work.

134 2 Sub-grid thickness fluxes

The mesoscale processes can be most strictly isolated with the help of an isopyc-
 nal model, as the adiabatic and diabatic processes are clearly distinguished in this frame-
 work. In this framework (stacked shallow water or isopycnal coordinate) the flow can be
 modeled using momentum and thickness equations (e.g. (Vallis, 2017; Loose, Marques,
 et al., 2023)). The thickness equation, the primary focus of our study, can be written
 for each layer as,

$$\partial_t h_n + \nabla \cdot (\mathbf{u}_n h_n) = 0, \quad (1)$$

135 where h_n and $\mathbf{u}_n = (u_n, v_n)$ are the thickness and velocity in the n^{th} layer respectively.

This equation can simulate the flow over the full spectrum of scales where its as-
 sumptions are valid, but with a finite grid size only a limited range of scales can be re-
 solved. Here we distinguish between scales that can be resolved and are too small to re-
 solved (sub-grid) using spatial filtering and coarse-graining $\overline{(\cdot)}$, as is routinely done in
 the large eddy simulation (LES) framework (Sagaut, 2005; Aluie et al., 2018). The high-
 pass signal after applying this spatial operation is referred to as *sub-grid* flows in this
 study. Consequently, the impact of sub-grid flow, on the resolved flows, can be elucidated
 in the thickness equation as,

$$\partial_t \overline{h_n} + \nabla \cdot (\overline{\mathbf{u}_n h_n}) = -\nabla \cdot (\overline{\mathbf{u}_n h_n} - \overline{\mathbf{u}_n} \overline{h_n}) \quad (2)$$

The impact of the sub-grid flow (e.g. $\mathbf{u}_n - \overline{\mathbf{u}_n}$) on the resolved flow (e.g. $\overline{h_n}$) arises as
 the divergence of a sub-grid flux on the RHS ($\nabla \cdot \mathbf{F}_n = \nabla \cdot (\overline{\mathbf{u}_n h_n} - \overline{\mathbf{u}_n} \overline{h_n})$). Hence-
 forth, we shall refer to

$$\mathbf{F}_n = \overline{\mathbf{u}_n h_n} - \overline{\mathbf{u}_n} \overline{h_n} \quad (3)$$

136 as the sub-grid scale (SGS) thickness flux, which will be the target of the parameteri-
 137 zations we develop below. It is common practice to represent this SGS thickness flux in
 138 terms of an eddy-driven stream function or velocity, as described in Appendix B. Also,
 139 in this theoretical framing we have assumed that our spatial filtering and coarse-graining
 140 commutes with the derivatives, which may not be true for all filter choices and near bound-
 141 aries (Moser et al., 2021); the exact choice of the operators used in this study is described
 142 in section 4.2.

143 Note that, considering resolved and sub-grid flows to the momentum equation would
 144 result in complementary SGS forcing terms in the momentum equation as well. How-
 145 ever in this study, we focus our attention only on the SGS thickness fluxes, as these cor-
 146 respond to one of the major parameterizations - the GM parameterization - in ocean mod-
 147 els. Data-drive parameterizations of SGS momentum forcing were considered recently

148 (Perezhogin et al., 2024; Zhang et al., 2023; Zanna & Bolton, 2020), and in follow-up work
 149 we plan to target both thickness and momentum SGS parameterizations simultaneously.

150 3 Machine learning model design and implementation

151 The goal of this work is to develop a data-driven parameterization for the SGS thick-
 152 ness fluxes in terms of the resolved variables. Here, we use a multi-layer perceptron (MLP)
 153 with a few hidden layers to approximate this functional relationship. These MLPs will
 154 be trained by using data from high-resolution (HR) simulations, which have been filtered
 155 and coarse-grained to diagnose the SGS thickness fluxes and the corresponding resolv-
 156 able fields. The machine learning parameterizations will be tested in both offline and on-
 157 line evaluation settings.

158 3.1 Parameterization function design

159 MLPs are universal function approximators, and can represent a vast space of func-
 160 tions. While this flexibility is powerful, it also increases the risk of overfitting, allowing
 161 the MLP to approximate data using functions that do not generalize. To avoid overfit-
 162 ting and allow for a degree of generalizability, we implement certain design constraints
 163 into the MLP.

Input features: Here we will search for functions of the form,

$$\mathbf{F}_n = f_\theta(\nabla\bar{\mathbf{u}}_n, \nabla\bar{h}_n, \Delta), \quad (4)$$

164 where $\nabla\bar{\mathbf{u}}_n$ is the velocity gradient tensor and $\nabla\bar{h}_n$ is the thickness gradient, both for
 165 the resolved fields. Δ is a measure of the grid scale, which allows the parameterization
 166 to be scale-aware. $f_\theta(\cdot)$ is a MLP function with unknown parameters θ , which need to
 167 be learned, that represents the two components of the SGS flux vector. Note that in the
 168 offline setting the inputs to this function will be the filtered and coarse-grained fields,
 169 while in the online setting the inputs will be resolved fields from the coarse-resolution
 170 simulation.

171 We found that functions of the above form can be trained to get remarkable offline
 172 success, but struggle when testing offline on data with distributional shifts (not shown).
 173 Adding some additional constraints discussed below, allows the model to generalize bet-
 174 ter to many more scenarios offline (Beucler et al., 2021).

175 **Lateral non-locality:** The GM parameterization and the velocity gradient model
 176 (VGM) parameterization, a common model used in the LES literature, are horizontally
 177 local (see Appendix C and Appendix D), i.e. parameterization output depends only on
 178 inputs from the same horizontal grid box (i, j) . Here, we relax this and allow for a small
 179 degree of non-locality in the horizontal, considering input information from regions sur-
 180 rounding the point where the prediction needs to be made. This is mathematically den-
 181 oted as $\mathbf{F}_{n,(i,j)} = f_\theta(\nabla\bar{\mathbf{u}}_{n,(I,J)}, \nabla\bar{h}_{n,(I,J)}, \Delta_{i,j})$, where $I = i + p$ and $J = j + q$, and
 182 p, q are integers in the range $(-m, m)$ – thus I, J correspond to the wider stencil around
 183 i, j . For a purely local model (1×1 stencil) $m = 0$, for a model with a 3×3 stencil
 184 $m = 1$, for a model with a 5×5 stencil $m = 2$ and so on. Note, that for all input stencil
 185 sizes, the prediction is always made only at the central point (i, j) . In principle, ver-
 186 tically non-local models can also be formulated, but these will not be considered here.

187 **Flow dependent coordinates:** We do not expect the sub-grid impacts to be co-
 188 ordinate dependent. However, when learning from data that comes from limited setups,
 189 a data-driven model may erroneously learn details about the coordinate. For example,
 190 if learning from data that comes from a f-plane channel simulation oriented in the x-direction,
 191 a data-driven model has the potential to learn that the SGS flux directed in the y-direction
 192 results in an available potential energy (APE) reduction. This model may fail if the setup
 193 is rotated by 90 degrees - even though we expect the impacts to not have changed.

To ward against this issue, we work in a flow-dependent, rather than a coordinate dependent, frame (Prakash et al., 2022). In particular, we rotate all our variables into a frame of reference oriented with the thickness gradients in each layer (see Appendix E for details). We denote variables in the flow dependent frame as (\cdot) . When working with laterally non-local inputs, the rotation is done with respect to the thickness gradient at the center point of the stencil. Also in this frame of reference, we implicitly reduce the number of inputs by one, as only the magnitude of the thickness gradient at the center of the stencil is now needed to quantify $\nabla\bar{h}_n$. This frame of reference is also conceptually advantageous, as the projection of the SGS thickness flux in the direction of the thickness gradient is responsible for the dissipation of resolved thickness variance, which is linked to the mean APE dissipation (Loose, Bachman, et al., 2023).

Non-dimensionalization and range limitation: The physics and sub-grid parameterizations should be invariant to the units of measurement. However, ML models are not unit invariant by default, and this property needs to be built in. Here this is achieved by casting all inputs and outputs into non-dimensional forms. In particular, we use the following non-dimensional forms: $\frac{\mathbf{F}_n}{\Delta^2|\nabla\bar{\mathbf{u}}_n||\nabla\bar{h}_n|}$, $\frac{\nabla\bar{\mathbf{u}}_n}{|\nabla\bar{\mathbf{u}}_n|}$, and $\frac{\nabla\bar{h}_n}{|\nabla\bar{h}_n|}$. Here $|\nabla\bar{\mathbf{u}}_n|$ is the For-

benius norm of the velocity gradient tensor; this is $\sqrt{\partial_x\bar{u}_n^2 + \partial_y\bar{u}_n^2 + \partial_x\bar{v}_n^2 + \partial_y\bar{v}_n^2}$ for a 1×1 stencil. For a larger stencil this takes the form

$$\sqrt{\sum_{p=-m}^{p=m} \sum_{q=-m}^{q=m} (\partial_x\bar{u}_{n,(i+p,j+q)}^2 + \partial_y\bar{u}_{n,(i+p,j+q)}^2 + \partial_x\bar{v}_{n,(i+p,j+q)}^2 + \partial_y\bar{v}_{n,(i+p,j+q)}^2)}.$$

Using the non-dimensionalization of the inputs using the magnitudes also ensures that all normalized input variables are limited in range between -1 and 1, helping constrain the input domain of the samples. There can still be gaps inside this multi-dimensional unit-sphere that were not sampled in the input data, but range limiting is still better than having unconstrained inputs. While not explicitly apparent, normalizing by the norm reduces the degree of freedom in input variable group by one. Apart from providing unit-invariance and range limiting, non-dimensionalization can potentially also provide generalization across some regimes where the energy levels are different but the underlying dynamics are similar.

Final functional form: The above considerations, result in the following :

$$\tilde{\mathbf{F}}_{n,(i,j)} = \Delta_{i,j}^2 |\nabla\bar{\mathbf{u}}_{n,(I,J)}| |\nabla\bar{h}_{n,(I,J)}| f_\theta \left(\frac{\widetilde{\nabla\bar{\mathbf{u}}}_{n,(I,J)}}{|\nabla\bar{\mathbf{u}}_{n,(I,J)}|}, \frac{\widetilde{\nabla\bar{h}}_{n,(I,J)}}{|\nabla\bar{h}_{n,(I,J)}|} \right), \quad (5)$$

where subscripts i, j and I, J have been included to make the non-locality of the model explicitly clear. No rotation is needed for the norms of the inputs, as the norm is invariant to coordinate rotation. While the above function choice is more restrictive than equation 4, it was chosen after some trial and error and we found that the functions estimated under these constraints are skillful.

3.2 Neural network architecture, hyper-parameters, and software

As mentioned above, we use a MLP to learn the function $f_\theta(\cdot)$. In this architecture the input layer is linked to the output layer through N_H hidden layers. Each hidden layer can have a different width (W_s , where $s \in (1, N_H)$), such that there are $\sum_{s=1}^{N_H} W_s$ hidden nodes. Each node applies a linear transformation to its inputs, followed by a non-linear activation function, which was chosen to be ReLU. The outputs of one layer serve as inputs to the next, enabling hierarchical feature learning. The final layer produces the output through a linear transformation.

We conducted a comprehensive sensitivity study on various MLP design and training choices, as detailed in Appendix F. Among all hyperparameters tested, we found model skill to be most sensitive to the total number of trainable parameters. For a given sten-

239 cil size and input/output configuration, performance improved with model size up to a
 240 threshold, beyond which additional parameters did not lead to further gains. Consequently,
 241 for discussion purposes in the results section, we restrict attention to models with ap-
 242 proximately the minimum number of parameters required to achieve maximal skill for
 243 each stencil size. We evaluate three MLP models that differ only in stencil size: 1X1, 3X3,
 244 and 5X5. Each model uses two hidden layers with 48 nodes per layer. As stencil size in-
 245 creases, the number of input features and thus the number of trainable parameters in-
 246 creases: from 2,786 (1X1), to 5,090 (3X3), to 9,698 (5X5). The models take non-dimensionalized
 247 velocity and thickness gradients as input and predict non-dimensionalized SGS thick-
 248 ness fluxes as output. Both inputs and outputs are rotated into local thickness gradi-
 249 ent coordinates. Also, in addition to the non-dimensionalization, all input and output
 250 features were also normalized by the order of magnitude of their standard deviations.

251 The details of the Double Gyre (DG) and Phillips 2 Layer (P2L) HR simulations
 252 are presented in the Section 4, and the details of the training set choices are the follow-
 253 ing. Training was performed on the first 2,048 snapshots from both the DG and P2L sim-
 254 ulations, sampled every 10 days and including spin-up. For each snapshot, data from mul-
 255 tiple filter and coarsening scales (details presented later) were used simultaneously. To
 256 give equal weight to each scale during training, data from finer filters were sub-sampled
 257 to align with the grid points of the coarsest filter. Offline evaluation was conducted on
 258 snapshots 2,4000 to 3,6000.

259 We used Python and JAX (<https://docs.jax.dev/>) for all our machine learn-
 260 ing pipelines. Specifically, we used Flax-Linen library ([https://flax-linen.readthedocs](https://flax-linen.readthedocs.io/)
 261 [.io/](https://flax-linen.readthedocs.io/)) for the design of our MLPs and used the Optax library ([https://optax.readthedocs](https://optax.readthedocs.io/)
 262 [.io/](https://optax.readthedocs.io/)) for optimization. The loss function was the mean absolute error in the non-dimensionalized
 263 outputs. Models were trained using the Adam optimizer with a learning rate of 0.01, and
 264 training was stopped when the validation loss failed to improve by more than 0.1% over
 265 10 consecutive epochs.

266 3.3 Implementation in MOM6

267 The MLP-based SGS thickness flux parameterization was implemented in MOM6
 268 via two new modules: an MLP module and a thickness flux prediction module. The MLP
 269 module, as the name suggests, reads a NetCDF file containing the model architecture,
 270 trained weights, and normalization factors, and performs inference as would be done by
 271 a standard feedforward MLP. This module is general-purpose and can be called from any-
 272 where within MOM6, enabling the integration of multiple MLP-based data-driven mod-
 273 els into the codebase.

274 The thickness flux prediction module incorporates the design choices described in
 275 Section 3.1. To keep the implementation simple and in light of the limited guidance in
 276 the literature regarding appropriate numerical schemes for such models we interpolate
 277 the necessary input fields to the centers of grid cells. For a 3X3 model, this means that
 278 each input in the 3X3 stencil surrounding the prediction point is evaluated at the 3X3
 279 grid cell centers. After predicting the two components of the thickness flux at the cen-
 280 tral point, the fluxes are then interpolated to the appropriate edge locations. We also
 281 introduced a non-dimensional coefficient (C_{ANN}) which can be used to adjust the strength
 282 of the parameterized flux if needed.

283 Fluxes at solid boundaries are set to zero to ensure volume conservation. We also
 284 observed that regions with very thin fluid layers could produce numerical artifacts. To
 285 suppress these, we modified the computation of the thickness gradient magnitude that
 286 multiplies the MLP prediction in Equation (5). Specifically, we replaced $|\nabla h_n|$ with $\left| h_n^2 \nabla \left(\frac{1}{h_n + \epsilon} \right) \right|$,
 287 which preserves the overall scaling in well-resolved regions but naturally drives the flux

288 toward zero in layers thinner than ϵ . This adjustment maintains the desired magnitude
 289 in most regions while ensuring numerical stability in thin layers.

290 While the P2L simulation does not contain vanishing layers, the lower layer in the
 291 DG simulation can vanish, as discussed in the next section. We tested the sensitivity of
 292 the results to different values of ϵ between 1 to 20 m and found the simulation outcome
 293 to be relatively insensitive to this choice.

294 While, in principle, the boundary conditions described above should suffice for lay-
 295 ered models (Killworth, 2001), MOM6 additionally enforces the constraint that the barotropic
 296 sum of the SGS thickness fluxes within a water column must be zero. This is accomplished
 297 by expressing the SGS fluxes in terms of a streamfunction (see Appendix Appendix B)
 298 and setting the streamfunction to zero at the boundaries using sophisticated tapering
 299 techniques (Ferrari et al., 2008, 2010). In our case, the situation is simpler, as we con-
 300 sider only two layers in the simulations evaluated in this study. Accordingly, we predict
 301 the lower-layer flux and satisfy this constraint by setting the upper-layer flux to be equal
 302 and opposite to the lower-layer flux. We also tested the parameterization without this
 303 constraint; while those simulations remained numerically stable, they frequently exhib-
 304 ited noisy solutions. Although enforcing a zero barotropic component has minimal im-
 305 pact on the overall energetics of the simulation (see appendix G4), it could play a role
 306 in lateral tracer transport. This effect is not studied here.

307 4 Data

308 4.1 Ocean Model Simulations

309 In this study, we work with two different idealized simulations in MOM6: Phillips
 310 2 Layer (P2L) - a 2 layer model of Phillips baroclinic instability, described in Hallberg
 311 (2013), and Double Gyre (DG) - 2 layer wind driven double gyre, described in Zhang et
 312 al. (2023). Both these setups have the minimum ingredients needed for the development
 313 of a rich baroclinic mesoscale field, while having a very simple vertical structure (Fig-
 314 ure 1). Some physical characteristics of the simulations are described in section 4.3 be-
 315 low. Simulations using both these setups were run over a range of resolutions. The high-
 316 est resolution output was filtered and coarsened for generating data to train and eval-
 317 uate the ML model offline, while the lower resolution simulations were used during on-
 318 line evaluation.

319 4.2 Processing of HR simulation data for ML training and evaluation

320 **Filtering and Coarsening:** To diagnose the input and output fields we processed
 321 the data using both filtering and coarse-graining. First for simplicity, we interpolated
 322 all the simulation prognostic variables onto the grid centers. The layer thickness and in-
 323 terface are already computed on the grid center, and u and v velocity components were
 324 linearly interpolated using the xgcm package (<https://xgcm.readthedocs.io/>). Also
 325 regions with layer thickness smaller than 20 m were masked and treated as land, which
 326 only impacted the lower layer in the DG simulation where the layer thickness vanishes
 327 at incropping locations. These centered and masked data were then filtered using a Gaus-
 328 sian filter, using the gcm-filters package (Loose et al., 2022). Specifically, we used the
 329 ‘simple fixed factor filter’ in gcm-filters ([https://gcm-filters.readthedocs.io/en/
 330 latest/examples/example_filter_types.html#simple-fixed-factor-filter](https://gcm-filters.readthedocs.io/en/latest/examples/example_filter_types.html#simple-fixed-factor-filter)), which
 331 does an area weighting. In the P2L simulation, with a Cartesian grid of size $\Delta_g = 4$ km,
 332 we used filter scales $L_f = 48, 100, 200$ and 400 km, where fixed filter factors of 12, 25,
 333 50 and 100 are used. In the DG simulation, with a non-uniform grid of size $\Delta_g = 1/20^\circ$
 334 this leads to variable filter scales of sizes $L_f = 0.55^\circ, 1.10^\circ, 2.20^\circ$ and 4.40° , which are
 335 approximately equal to the scales used for filtering P2L.

336 After all the filtered variables and the corresponding SGS fluxes were computed,
 337 the data was further coarsened using box averages. This step only leads to a data re-
 338 duction, and for convenience we do not account for the SGS fluxes resulting from this
 339 operation. This is sufficient since the fields were already filtered, and the additional fluxes
 340 resulting from this coarsening operation are found to be very small (not shown). We de-
 341 cided to choose the ratio between the filtering scale and the coarsening scale, to be 5.
 342 For example, a filter scale $L_f = 200$ km was combined with a coarsening scale of $\Delta_c =$
 343 40 km. This choice of the filter to grid ratio ($FGR=L_f/\Delta_c$) is based partially on con-
 344 sidering the spectra of ocean models, which often show a damped energy level or numer-
 345 ical artifacts emerging at scales larger than the grid scale - often upto 5 times in size (e.g.
 346 notice the small-scale bump in the EKE spectrum in Figure 7). This is a heuristic choice,
 347 and can be explored in more detail in future work. For the rest of this study we indi-
 348 cated the four filter scales nominally with $L_f = 50, 100, 200,$ and 400 km, and correspond-
 349 ing coarse grid scales nominally with $\Delta_c = 10, 20, 40$ and 80 km. However this is only
 350 for the convenience of presentation, and in the computations, where Δ_c is needed (e.g.
 351 equation 5), the actual coarse grid scales were used.

352 Currently there is no prescribed way to process data from a HR simulation to make
 353 it match a LR simulation in some objective sense. Hence, our approach to diagnosing
 354 filtered and coarsened data is relatively adhoc, and based on pragmatism. We believe
 355 and hope that the impact of these choices would likely be minimal, with the acknowl-
 356 edgement that a big shortcomings of these ML models in online settings would likely arise
 357 from the fact that the distribution of a low resolution simulation would be shifted rel-
 358 ative to a simplistically filtered version of a high resolution simulation, and some degree
 359 of tuning may be required to address this. Also, in future work more care can be taken
 360 for dealing with staggered grids, precisely accounting for the separate contributions from
 361 filtering and coarsening operations, and for handling boundaries differently when the fil-
 362 ters do not commute with the gradients.

363 **Layer thickness decomposition:** When computing the thickness fluxes, special
 364 care was taken when the bottom topography (η_b) was not flat (in the DG case). We did
 365 not want to filter the bottom topography when filtering thickness, since the topography
 366 present in a coarse model is not a filtered version of the high resolution topography. So
 367 we chose to use the condition that $\bar{\eta}_b = \eta_b$. Thus, thickness fluxes were dealt with by
 368 filtering interface heights only. To be more precise, in the bottom layer the filtered thick-
 369 ness would be $\bar{h}_N = \bar{\eta}_{N-1/2} - \eta_b$ and the filtered advection would be $\bar{\mathbf{u}}h_N = \bar{\mathbf{u}}\bar{\eta}_{N-1/2} -$
 370 $\bar{\mathbf{u}}\eta_b$, where $\eta_{N-1/2}$ is the upper interface height of the bottom layer N .

371 Further, we decomposed thickness gradients into a steady and a deformable parts
 372 ($\nabla\bar{h}_n = \bar{h}_N^{deformable} + \nabla\bar{h}_n^{steady}$). When applied to a layer where the lower interface is
 373 topography, (e.g. $\bar{h}_N = \bar{\eta}_{N-1/2} - \eta_b$) we get $\nabla\bar{h}_N^{steady} = -\nabla\eta_b$ and $\nabla\bar{h}_N^{deformable} =$
 374 $\nabla\bar{\eta}_{N-1/2}$. In other layers, the steady contribution is zero ($\nabla\bar{h}_n^{steady} = 0$) and the full
 375 layer thickness contributes to the deformable part ($\nabla\bar{h}_n^{deformable} = \nabla\bar{\eta}_{n-1/2} - \nabla\bar{\eta}_{n+1/2} =$
 376 $\nabla\bar{h}_n$). This decomposition allows us to distinguish the impact of dynamic layer thick-
 377 ness variations and bottom topography on the SGS fluxes. In this study, we only use the
 378 deformable contribution as inputs to our MLP, and henceforth replace the notation for
 379 the deformable part $\nabla\bar{h}_N^{deformable}$ by $\nabla\bar{h}_N$ for simplicity in most places, unless other-
 380 wise noted.

381 4.3 Physical characteristics of the simulations and the filtered data

382 4.3.1 Data distributions

383 Both the P2L and DG simulations produce a turbulent flow field, with a rich ar-
 384 ray of eddies and jets (Figure 1). The magnitude of the SGS fluxes are greater in the
 385 top layer than the bottom layer, and are larger in P2L than DG (Figure 2). Also the SGS

386 fluxes vary by three to four orders of magnitude in each layer, and the magnitude of these
 387 fluxes increase by about an order of magnitude from filter scales of 50 km to 400 km (coarse
 388 grid scales of 10 km to 80 km). In contrast, the non-dimensionalized fluxes have a much
 389 narrower distribution, with little variation of the median across filter scales and layers
 390 – indicating that the non-dimensionalization choices made here are quite successful at
 391 collapsing the data distribution and may help with generalization. However, note that
 392 the width of non-dimensionalized flux distribution does increase slightly with filter scale,
 393 indicating that there may still be some room for further improved non-dimensionalization
 394 factors in the future, which may be able to address this scale dependence.

395 Parameterizations in Z-level models and MOM6 impose that the eddy driven stream
 396 function takes a boundary condition of the zero at the surface, which is equivalent to say-
 397 ing that there is no barotropic (depth integrated) SGS flux (see Appendix B). This condi-
 398 tion is not naturally satisfied by the diagnosed data (red distributions in Figure 2), and
 399 is not even expected based on the eddy-mean decomposition of the thickness equation
 400 in layered models (Killworth, 2001). In the diagnosed data, we notice that the fluxes in
 401 the two layer have a very slight opposing tendency, which slightly increases with larger
 402 filter scales. Only in the long time average, and if the mean flow is weak, would we ex-
 403 pect the vertical sum of the eddy thickness fluxes to go to zero based on volume conser-
 404 vation. However, it is worth noting that even though the depth integrated SGS thick-
 405 ness flux is far from negligible, its impact on the APE tendency, which is the primary
 406 target of our parameterization, is very small (see Appendix G4).

407 The distributions of other model fields, particularly those relevant for our neural
 408 network design, are shown in Figure 3. The velocity gradients have a large range across
 409 scales, and their magnitude decreases with increasing filter scale. Generally, the upper
 410 layer has stronger velocity gradients than lower layer, and the P2L simulation has stronger
 411 velocity gradients than the DG simulation. In contrast to velocity gradients, the deformable
 412 thickness gradients vary less with filter scale. Additionally, since the surface variations
 413 are much weaker than the interface variations (Figure 3c), the deformable thickness gra-
 414 dients are essentially the same for the two layers. Also, the deformable thickness gra-
 415 dients are slightly weaker in the DG simulation relative to the P2L simulation. The bot-
 416 tom topography slopes in the DG simulation are very strong relative to the interface vari-
 417 ations, which was one reason for us to decompose the thickness gradients into its steady
 418 and deformable contributions.

419 **4.3.2 Bulk properties**

420 The rich turbulent eddies impact the mean or large spatial and temporal scale flow,
 421 and the target of traditional parameterizations has been to skillfully model some of these
 422 effects. Here we describe what some of these feedback are when the eddies are resolved.

423 In the P2L simulation the mean state is a zonal jet, which is sustained by slowly
 424 relaxing the middle interface to a sloping state and thus the relaxation works as a source
 425 of APE (Hallberg, 2013). The eddies in this simulation work to flatten this interface, thus
 426 removing the APE that is input by the relaxation forcing. To achieve this the eddies flux
 427 volume to the north in the upper layer and to the south in the lower layer, generating
 428 an eddy driven overturning circulation. This overturning is sustained because the relax-
 429 ation also leads to a diapycnal transformation of watermasses from one layer to the other.
 430 When the eddies are not resolved or only partially resolved, this eddy-driven overturn-
 431 ing circulation weakens and a parameterization is needed to ensure that the appropri-
 432 ate levels of APE are removed and the overturning circulation is maintained (Hallberg,
 433 2013).

434 In the DG simulation (Zhang et al., 2023), the mean state is maintained by forc-
 435 ing with a steady wind stress. The wind stress peaks at the intermediate latitude (40N)
 436 and generates a region of Ekman downwelling to the south and Ekman upwelling to the

437 north, which pushes and pulls the middle interface ("thermocline") to create an inter-
 438 face slope. The sloping interface has APE that is siphoned out by the mesoscale eddies,
 439 which work to reduce the APE and thus counteract the impact of the wind.

440 5 Results

441 Data-driven models may be evaluated along two interconnected aspects: offline skill
 442 and online skill. Offline skill refers to the accuracy in predicting SGS thickness fluxes
 443 for a given set of inputs, with the reference or "truth" diagnosed from high-resolution
 444 simulations. In contrast, online skill assesses the ability of a data-driven parameteriza-
 445 tion to improve the fidelity of a lower-resolution simulation, where the truth is defined
 446 in terms of large-scale behavior - either from high-resolution models or observations.

447 Here, we first evaluate the offline skill of the machine learning model on filtered and
 448 coarsened data. We then turn to the arguably more important question: how well does
 449 the model perform online, when it is embedded within a simulation and actively inter-
 450 acts with and modifies the resolved state?

451 5.1 Offline Evaluation

452 Here we discuss three particular MLP models that differ only in stencil size: 1X1,
 453 3X3, and 5X5 (details in section 3.2). We also contrast these models against the GM pa-
 454 rameterization and the VGM parameterization, where the free parameters in these con-
 455 ventional schemes were estimated using least-squares fitting to the SGS thickness fluxes
 456 for each parameterization separately at each filter scale and for each simulation setup.
 457 This is in contrast to the MLPs, which were trained across the entire range of filter scales
 458 and simulations simultaneously. Given the large variation in the estimated parameters
 459 for the conventional parameterizations, their performance would be worse if they were
 460 trained across the entire range of data simultaneously.

461 **Point-wise skill:** Mostly, all the MLP models demonstrate relatively high skill
 462 in predicting the SGS thickness fluxes point-wise. As an example, the true and predicted
 463 SGS fluxes from the 3X3 model at a filter scale of 100 km are shown in Figure 4. The
 464 MLP model does extremely well, and produces the right patterns and magnitudes of both
 465 the along and across thickness gradient components of the SGS flux in both layers. Note
 466 in the figure that the error is so small that it had to be multiplied by 5 to bring it to the
 467 same color scale as the SGS flux.

468 In a more quantitative sense, we find that the MLP skill, compared across all MLP
 469 configurations considered here and quantified using R2 or correlation score (Appendix
 470 A), depends on the input stencil size and the filter scale (Figure 5); the skill also depends
 471 on other factors as discussed in Appendix G but those sensitivities can be alleviated with
 472 enough data or free parameters. The ML model skill increases with stencil size, with a
 473 large improvement in going from 1X1 to 3X3 and a smaller improvement when going fur-
 474 ther to 5X5. Also, model skill decays as filter scales get larger, which is more rapid in
 475 the case of DG relative to P2L. We think that this might have to do with the different
 476 deformation radii, which is on average 40 km in P2L and 20 km in DG (Figure H1), and
 477 ML models may be less skillful as the filter scale gets much larger than the deformation
 478 radius or larger than the largest resolved eddies.

479 In figure 5, we also contrasted the skill of the MLPs against the conventional pa-
 480 rameterizations. The VGM parameterization, which is sometimes used in LES studies,
 481 has essentially the same point-wise skill as our 1X1 ML model. The GM parameteriza-
 482 tion in contrast has very little point-wise skill (both R2 and correlation are usually less
 483 than 0.2), which only marginally increases at larger filter scales. This is to be expected,
 484 as the GM parameterization is a bulk model and not designed to have skill in produc-

485 ing the right local structural patterns in the SGS fluxes. Notice that in contrast, the skill
 486 of the 3X3 and 5X5 models is significantly higher than both the 1X1 ML model and the
 487 VGM parameterization.

488 **Bulk (time-averaged) skill:** While pointwise skill is a useful metric, we are of-
 489 ten (also) interested in ensuring that our parameterizations produce the appropriate bulk
 490 effects. One way to quantify the bulk effects is in terms of time averages. In Figure 6
 491 we consider the skill over different temporal averaging windows, here the skill is quan-
 492 tified for both layers and flux components.

493 We expect that the GM parameterization skill may improve in this bulk sense as
 494 the time averaging duration is increased, and as expected this is found to be the case.
 495 In Figure 6 last column, this effect is seen quite clearly for the P2L data but not for the
 496 DG data. Even for the P2L case the correlation skill only rises to 0.5, rather than 1, which
 497 is because GM only predicts the along gradient fluxes, and so when averaging skill over
 498 along and across gradient fluxes, we can only achieve a maximum of 0.5. The reason for
 499 the discrepancy between DG and P2L arises because GM did not turn out to be a good
 500 model for upper layer fluxes in the DG case (even when quantified just in terms of cor-
 501 relation), which might be a result of mean flows and inhomogeneity in the turbulent statis-
 502 tics. The skill of the GM model on the lower layer along gradient fluxes in the DG data
 503 is higher, but still not as large as for P2L data (not shown).

504 The 1X1 MLP, and similarly the VGM parameterization (Khani & Dawson, 2023),
 505 have R2 skill decrease with increasing temporal averaging. In contrast, the correlation
 506 is not impacted, suggesting that the decrease in skill has less to do with the functional
 507 form and more to do with the parameterization coefficient or amplitude. In contrast the
 508 3X3 and 5X5 MLPs have almost no impact on either R2 or correlation skill with aver-
 509 aging, if anything there is a very slight increase in skill at longer temporal averaging. The
 510 fact that the skill score is usually close to 1, except for larger filter sales in DG, also shows
 511 that these models do very well in predicting both the along and across fluxes.

512 Overall, we found that the offline skill of the MLP models, particularly those with
 513 wider stencils, is very promising. Also, we show in Appendix Figure G4 that a MLP trained
 514 on data from one simulation shows relatively high skill when tested on the unseen simu-
 515 lation. Thus, these new models are scale and context aware, and no retraining or tun-
 516 ing is needed when testing over different datasets. This is contrast to the traditional mod-
 517 els, which are unable to match the MLP skill even after the corresponding coefficients
 518 were estimated separately for each scale and dataset. Thus, we are compelled to eval-
 519 uate the performance of these MLPs in an online setting.

520 In the online setting we only evaluate the skill of the 3X3 MLP discussed above,
 521 as this model provides a good compromise between computational cost and offline skill
 522 (skill of the 5X5 MLP is only marginally better than the 3X3 MLP). Also since the of-
 523 fline performance of the 1X1 model (and VGM) is worse and degrades in a bulk sense,
 524 we chose to not evaluate it either.

525 5.2 Online Evaluation - Phillips 2 Layer

526 We first test the MLP in the P2L simulation, which is the simpler of the two sim-
 527 ulation setups considered in this study. To assess the sensitivity of this setup to model
 528 resolution and parameterizations, we performed a suite of simulations (see Appendix H).

529 In this section, we focus on the 20 km simulations – both parameterized and un-
 530 parameterized – and compare them with a 4 km HR simulation. The deformation radii
 531 in this setup range from 25 to 50 km, placing the 20 km grid in the “gray zone” where
 532 mesoscale eddies are only partially resolved, whereas the 4 km grid resolves them more
 533 fully. The parameterized simulations were selected to approximately match the total over-

534 turning transport of the HR case. For the GM parameterization, we include two con-
 535 figurations: one with low diffusivity ($1000 \text{ m}^2/\text{s}$), often cited as a canonical value in the
 536 literature, and another with high diffusivity ($8000 \text{ m}^2/\text{s}$), chosen because it yields an over-
 537 turning close to the HR simulation. These values are not fine-tuned to exactly match
 538 the HR overturning, but instead represent two qualitatively distinct regimes. Notably,
 539 $1000 \text{ m}^2/\text{s}$ also marks a threshold beyond which the GM scheme begins to strongly damp
 540 the resolved eddy field (also see Appendix H). In contrast, the MLP-based parameter-
 541 ization with a tuning coefficient (C_{ANN}) of 1 produced overturning transport closely aligned
 542 with the HR simulation. All simulations reached a spun-up state within 2 years, and time-
 543 averaged statistics are computed over years 2 to 10.

544 Snapshots of upper layer relative vorticity and the EKE spectrum averaged tempo-
 545 rally and meridionally in all simulations evaluated in this section are shown in Fig-
 546 ure 7. The variability in the 20 km simulations without any parameterization and with
 547 the MLP is very similar, while the addition of the GM parameterization leads to a sub-
 548 stantial reduction in flow variability. The low GM diffusivity simulation permits some
 549 eddies, while the eddies are entirely suppressed in the high GM diffusivity simulation.
 550 Also, the HR filtered and coarsened simulation state matches the low-resolution unpa-
 551 rameterized simulation at large scales, but has lower energy levels at smaller scales. This
 552 is the result of the specific properties of the Gaussian filter that was chosen to filter the
 553 simulation, and more refined filters could definitely be employed if needed.

554 The peak overturning transport in the upper layer is shown in Figure 8a (the lower
 555 layer's overturning is identical but with the opposite sign). The HR simulation produces
 556 approximately 13 Sv of transport, with 9.5 Sv attributed to scales larger than the filter
 557 scales and 3.5 Sv from SGS fluxes. In contrast, the unparameterized 20 km simulation
 558 produces only about 10.5 Sv of transport. The 20 km simulation with the MLP sustains
 559 about 12.5 Sv of transport, reducing the resolved component to 9.5 Sv while adding around
 560 3 Sv from parameterized fluxes. The low GM diffusivity case similarly reduces the re-
 561 solved transport marginally, but is unable to generate enough parameterized flux to match
 562 the total overturning transport of the high-resolution simulation. As the GM diffusiv-
 563 ity is increased, the resolved transport is drastically reduced, with only a slight increase
 564 in the total transport. When the GM diffusivity becomes large enough, the total trans-
 565 port can match that of the high-resolution simulation, but at the expense of completely
 566 eliminating the resolved contribution. Note that we can further fine tune the MLP and
 567 the high GM diffusivity case to exactly match the HR simulation. However, this is not
 568 possible for the low GM diffusivity case because the resolved transport drops much more
 569 rapidly than the rate at which the parameterized transport increases with changing the
 570 coefficient (can also see Figure 7 in Hallberg (2013)).

571 While the parameterizations were tuned to approximately match the overturning,
 572 here we focus on contrasting their impact on other relevant metrics (details of these met-
 573 rics are described in Appendix G). The kinetic energy (KE) and available potential en-
 574 ergy (APE) from various contributions are shown in Figures 8b and c. The unparam-
 575 eterized 20 km simulation exhibits lower KE and APE than the HR simulation, as ex-
 576 pected, and the LR simulation KE and APE are close to the filtered KE and APE from
 577 the filtered HR data. The 20 km simulation with the MLP results in a slight reduction
 578 in the EKE and EAPE, but the overall KE and APE is roughly in line with the KE and
 579 APE from the filtered HR simulation. The low GM diffusivity simulation shows lower
 580 KE and APE than both the unparameterized and MLP-based simulations, while the high
 581 GM diffusivity case has no EKE or EAPE. Since in this simulation setup the interface
 582 height is restored to a prescribed state, the MKE and MAPE remain nearly unchanged
 583 across the different setups. Only in the high GM diffusivity case does the parameteri-
 584 zation forcing become large enough to cause a very small reduction in MKE and MAPE,
 585 and the EKE and EAPE are completely wiped out.

586 The tendency of the APE arising from the SGS or parameterized fluxes is shown
 587 in Figure 8d. We display the contribution to both the MAPE tendency and the EAPE
 588 tendency. In the filtered HR simulation, the SGS fluxes contribute just over half of their
 589 tendency towards reducing the MAPE. The 20 km simulation with the MLP parame-
 590 terization produces an APE tendency that is relatively close to the HR filtered case, with
 591 a slightly larger impact on the MAPE. In contrast, the low GM diffusivity case leads to
 592 a disproportionately large impact on the EAPE, without sufficiently reducing the MAPE.
 593 The high GM diffusivity case produces a similar MAPE tendency to the low diffusivity
 594 case but has no impact on the EAPE, as no eddies remain to be damped out.

595 While the APE tendency analysis is illustrative, it has some limitations due to the
 596 changes in the simulation state across different cases. To further emphasize this point,
 597 we also evaluated the APE tendency that would arise from the MLP in the 20 km sim-
 598 ulation, where the MLP was not actually coupled to the resolved state of the simulation
 599 (see bar labeled 20 km in Figure 8d). In this scenario, the APE tendency is much larger
 600 than in the filtered case, likely because this unparameterized simulation has a higher eddy
 601 kinetic energy (EKE) than the filtered EKE. The non-linear interaction between the pa-
 602 rameterized and resolved flow makes it difficult to predict a priori how the system will
 603 respond to the parameterization. This non-linearity is also evident in the sensitivity study
 604 plots shown in Appendix H, where the response to the parameterization coefficients is
 605 non-monotonic. Note that at eddy-permitting resolutions, even the response to the GM
 606 parameterization is non-trivial and non-monotonic.

607 In summary, while both the MLP and GM parameterizations can be tuned to pro-
 608 duce approximately the correct overturning circulation in the P2L simulation, only the
 609 MLP is able to achieve this without significantly damaging the resolved flow and eddies.
 610 In contrast, as also shown by (Hallberg, 2013), the GM parameterization excessively dis-
 611 sipates the eddies.

612 5.3 Online Evaluation - Double Gyre

613 Next, we test the MLP parameterization in the DG simulation, a canonical system
 614 for studying wind-driven gyre dynamics. In contrast to the P2L setup discussed earlier,
 615 the DG exhibits strong boundary currents and pronounced spatial inhomogeneity in eddy
 616 statistics. Additionally, unlike the P2L system, the mean state here is not maintained
 617 through relaxation, but is rather a result of balance between winds and eddies. These
 618 differences lead to two key consequences that are different from P2L: (i) mean-state bi-
 619 ases can emerge as resolution and parameterizations are varied, and (ii) the system is
 620 purely adiabatic, with no overturning circulation.

621 Similar to the P2L case, we conducted a suite of simulations to assess the sensi-
 622 tivity of the DG setup to both resolution and parameterization coefficients (see Appendix
 623 H). Here, we focus on the $1/5^\circ$ (~ 20 km) simulations, which only marginally resolve the
 624 deformation radius - ranging between 5 and 30 km in this configuration. In addition to
 625 the unparameterized baseline, we analyze simulations that employ MLP and GM param-
 626 eterizations. The coefficients for these parameterized runs were selected to minimize the
 627 mean state error in sea surface height (SSH), which is strongly correlated with thermo-
 628 cline depth in this system.

629 The mean sea surface height (SSH) and kinetic energy (KE) fields, averaged over
 630 years 3-13 of the simulations, are shown in Figure 9. In the $1/5^\circ$ simulation without any
 631 parameterization, a standing eddy forms just downstream of the boundary current sep-
 632 aration point - a region that also exhibits stronger flow than in the high-resolution ref-
 633 erence simulation. This eddy feature vanishes, and the mean state bias is reduced, in both
 634 the MLP- and GM-parameterized simulations, since the parameterizations coefficients
 635 were explicitly chosen to reduce this bias. However, consistent with the P2L results, the

636 MLP achieves this correction without the substantial loss of KE observed in the GM sim-
637 ulation.

638 This result is quantified in the KE and APE metrics shown in Figure 10. The mean
639 state of the $1/5^\circ$ simulation without parameterization is overly energetic compared to
640 the HR filtered simulation, as evident in both the MKE and MAPE. While both param-
641 eterizations reduce this excess energy, the GM parameterization does so at the cost of
642 a much larger reduction in EKE and EAPE compared to the MLP parameterization.

643 To further assess the role of the parameterizations in influencing the mean flow and
644 eddies, we examine the APE tendency induced by SGS fluxes – both in the filtered HR
645 simulation and in the parameterized fluxes of the coarser simulations (Figure 11). The
646 APE tendencies affecting both the MAPE and EAPE exhibit qualitative similarities be-
647 tween the filtered HR and the MLP parameterized simulations, generally acting to re-
648 duce APE. However, this impact is not spatially or temporally uniform and in many in-
649 stances there is even APE gain, resulting in localized regions of APE gain even in the
650 10-year mean shown here. In contrast, the GM parameterization acts as a sign-definite
651 sink of APE, producing a much stronger and more widespread reduction in both MAPE
652 and EAPE compared to the MLP parameterization or the diagnosed tendencies from the
653 filtered HR simulation.

654 In summary, as in the P2L case, the MLP parameterization outperforms the GM
655 parameterization in the DG simulation. It effectively reduces the mean state bias with-
656 out causing an excessive suppression of eddy energy.

657 **6 Discussion and Conclusions**

658 In this work we developed and implemented a data-driven parameterization for sub-
659 grid scale (SGS) thickness fluxes produced by mesoscale eddies. This was achieved by
660 training a relatively small multi-layer perceptron (MLP) to learn a functional relation-
661 ship between the gradients of the large-scale/resolved fields and the SGS fluxes, using
662 data from high resolution (HR) simulations. By introducing features like lateral non-locality,
663 coordinate invariance, and non-dimensionalization into the MLP design, we were able
664 to produce a more generalizable and stable data-driven parameterization (Perezhogin
665 et al., submitted). Of these features, the lateral non-locality and the non-dimensionalization,
666 particularly the aspect that produces range-limited inputs, were found to be the most
667 important design choices. The trained models have very high offline skill (Figure 5), even
668 when testing on data coming from unseen simulations (Figure G4). The skill relatively
669 degrades at scales larger than the largest eddies, or as we transition from eddy permit-
670 ting to non-eddy resolutions, but even at these scales the offline skill is comparable
671 or higher than traditional approaches like an appropriately tuned Gent-McWilliams (GM)
672 parameterization (Figure 6).

673 This new data-driven parameterization was implemented into GFDL’s Modular Ocean
674 Model 6 (MOM6) and tested in two idealized simulation setups: Phillips 2 Layer (P2L)
675 and Double Gyre (DG), where baroclinic mesoscale eddies play a first order role in the
676 dynamics. In both these setups the MLP enhanced the simulation performance at coarse
677 resolutions (grid scales on the order of the deformation radius or coarser), reducing bi-
678 ases in aspects like the meridional overturning transport and mean state. In the MOM6
679 implementation, we also introduced a non-dimensional tuning parameter that controls
680 the global amplitude of the parameterization. Sensitivity studies showed that $O(1)$ val-
681 ues of this parameter, values of 1 for P2L and 0.5 for DG, were optimal in online set-
682 ting across all eddy permitting resolutions. At coarser resolutions the values needed to
683 be slightly adjusted, a value of 2 for P2L and 0.75 for DG seemed optimal. In contrast,
684 the GM diffusivity had to be adjusted for every resolution and setup individually, with

685 optimal values of the diffusivity ranging between three orders of magnitude (10^2 to 10^4
 686 m^2/s) across resolutions and setups.

687 The development of the GM parameterization more than three decades ago provided a step change in the quality and fidelity of ocean simulations. However, it has been
 688 clear since the beginning that while GM parameterization is phenomenologically appropriate, reduction of APE and adiabatic conservation of watermass volume are appropriate,
 689 reduction of APE and adiabatic conservation of watermass volume are appropriate bulk expectations for the effects of mesoscale eddies, it has structural shortcomings. Much of the work over the past decades has gone towards improved parameter estimation
 690 (Visbeck et al., 1997), but progress towards reducing the structural errors has been very limited (R. D. Smith & Gent, 2004). Our work attempts to address this gap, and
 691 provide a data-driven model that seems to have lower structural errors and provides a path towards improving the GM parameterization along a new axis.
 692
 693
 694
 695
 696

697 One of the well known drawbacks of using the GM parameterization is its tendency to dissipate the resolved eddies, reducing the effective resolution of a simulation (Hallberg,
 698 2013). This has led to development of adhoc fixes, like turning off the GM parameterization using a resolution function, or not using the GM parameterization at all in simulations
 699 that run in the gray zone resolution (Adcroft et al., 2019). Our parameterization does not suffer from this problem, and is able to improve properties of the large scale
 700 without being overly detrimental towards the resolved eddies. One major advantage of the GM parameterization is that it is essentially guaranteed to be stable, which can not
 701 be claimed unequivocally for our data-driven parameterization. All the simulations we tested were stable, but this may not hold to be true if the tuning coefficient is pushed
 702 to larger values.
 703
 704
 705
 706
 707

708 Our data-driven parameterization has produced very promising results, but a few aspects can likely be improved. Firstly, as is common with machine learning models, there
 709 is a spectral bias in the predictions. This means that the MLP’s offline performance degrades at scales where the signal variance is low (Figure G3 and G4, and also implicit
 710 in Figure 4). This did not end up being an issue for model stability in our online tests, likely because the thickness variance has a tendency to cascade down-scale and the errors
 711 get cascaded to model dissipation scales. Improved loss functions or experimenting with smoother non-linearities could help alleviate this issue. Secondly, it came almost
 712 as a surprise that the model performed so well across resolutions and setups with such limited class of inputs. We think that at eddy permitting resolutions this is the case because
 713 the structure of largest eddies that are resolved already contain a lot of information about dynamically important environmental conditions, which shape the mesoscale eddy field.
 714 Reduction in offline skill at scales larger than the size of the largest eddies and change in tuning parameters at non-eddy permitting resolutions led us to this hypothesis.
 715 This suggests that in future work it may be worth paying more attention to these scales, and potentially introducing more model inputs or vertical non-locality to be able to perform
 716 well over a larger range of scales without much tuning. Thirdly, our parameterization mainly acts on reducing the APE in the system, and is not designed to improve the MKE or EKE
 717 of the resolved state that may arise if the removed APE was appropriately cascade upscale in an inverse KE cascade. Coupling our parameterization with a backscatter parameterization
 718 in an energetically consistent manner could be investigated to produce further improvements in ocean models. Lastly, we chose to use a reasonably small MLP to keep the computational
 719 burden due to the parameterization low in principle. However, we have not performed a comprehensive testing and optimization to get the best possible computational performance for our model in MOM6, yet.
 720
 721
 722
 723
 724
 725
 726
 727
 728
 729
 730
 731
 732

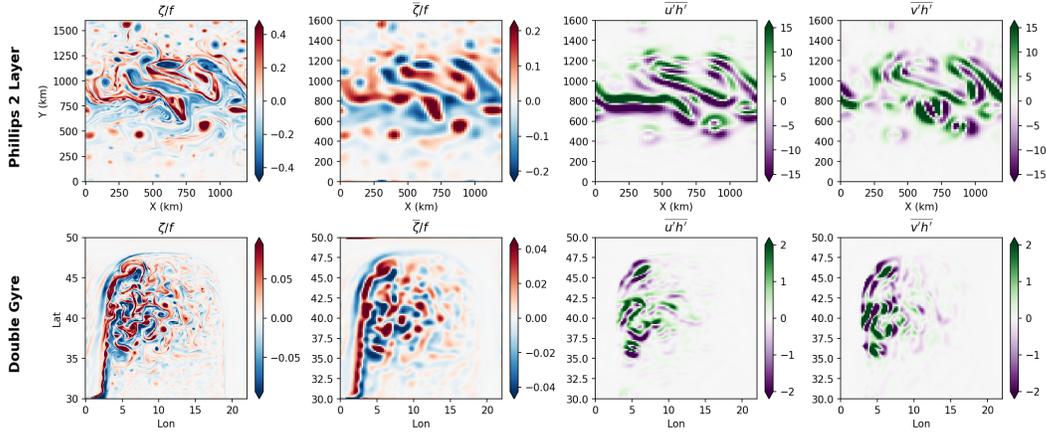


Figure 1. Snapshots exemplifying the high-resolution vorticity (left), filtered vorticity (second column), and sub-grid fluxes (last two columns) from the Phillips 2 layer (top) and Double Gyre (bottom). The filtered fields are shown from the case of coarse-graining scale of 20 km. Only the upper layer data is shown, and different panels have different color ranges.

733

Open Research

734

All the code to generate the simulation data and machine learning data, do the machine learning training, and do the analysis and figure generation can be found at https://github.com/dhrubwalwada/mesoscale_buoyancy_param_ML.

735

737

Acknowledgments

738

This project is supported by Schmidt Sciences, LLC. DB also wants to acknowledge Ryan Abernathy, now at EarthMover PBC, for providing him advice and guidance, and helping him develop his research interests in the area of mesoscale parameterizations for numerical models. This work and DB's career likely would not have shaped out the way they did without Ryan's constant support, encouragement, and enthusiasm.

741

743

Appendix A Offline Skill Metrics

To evaluate the skill of our ML model in an offline setting we use two main skill metrics. The first is the coefficient of determination or R2 skill, defined as,

$$R^2 = 1 - \frac{\overline{(y - \hat{y})^2}}{\overline{(y - \bar{y})^2}}, \quad (\text{A1})$$

744

y is the truth value and \hat{y} is the prediction. The $\overline{(\cdot)}$ corresponds to average over all samples being considered, which is chosen to be over both flux components from both layers, full spatial domain, and all temporal snapshots in the test data (unless indicated otherwise). The R2 skill will be 1 when prediction is perfect and reduces as prediction gets worse.

745

746

747

748

The second metric is the Pearson correlation coefficient,

$$C = \frac{\overline{(y - \bar{y})(\hat{y} - \bar{\hat{y}})}}{\sqrt{\overline{(y - \bar{y})^2} \overline{(\hat{y} - \bar{\hat{y}})^2}}}, \quad (\text{A2})$$

749

which is 1 when the truth and the prediction are perfectly correlated, -1 for inverse correlation, and 0 for no correlation.

750

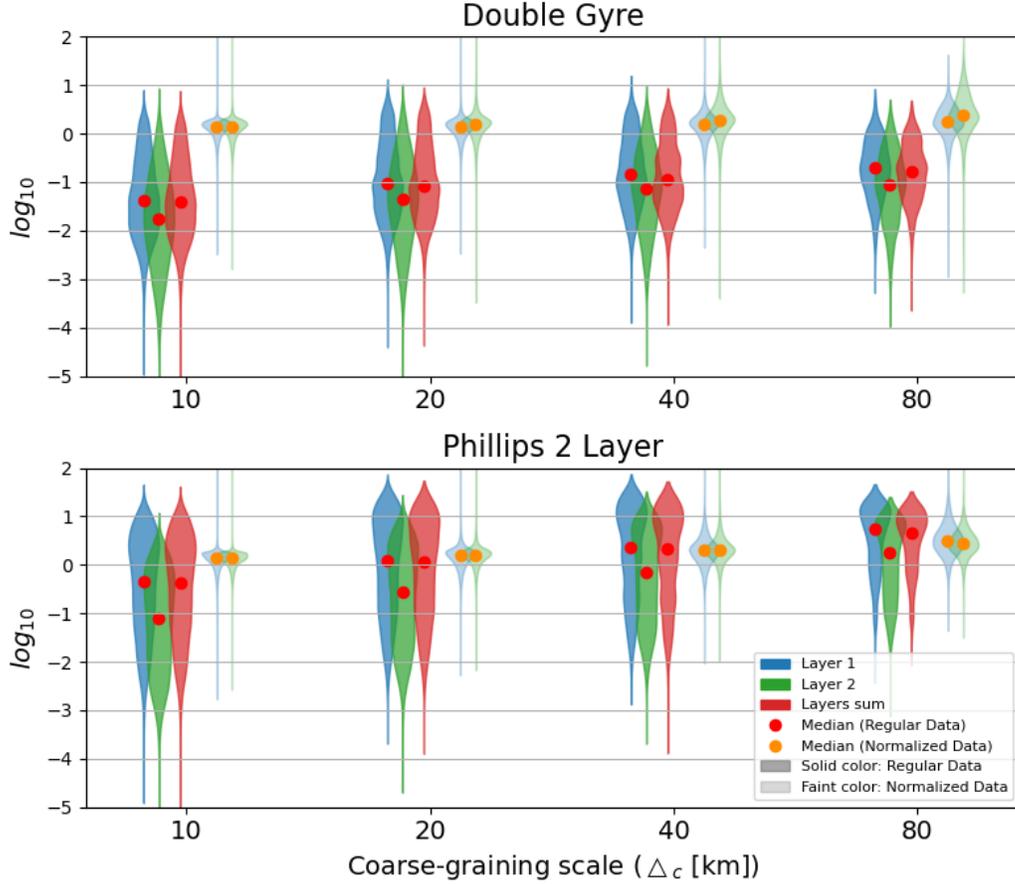


Figure 2. Distribution of output data: Distributions of the logarithm of the thickness flux magnitudes ($|\mathbf{F}_n|$) in both layers and their sum (barotropic contribution), and the normalized thickness flux ($\frac{|\mathbf{F}_n|}{\Delta_c^2 |\nabla \mathbf{u}_n| |\nabla h_n|}$) in both layers for the different experiments (top and bottom panel) and different coarse-graining scales (indicated on the x-axis). The legend in the lower panel is used to indicate the different elements in both the panels. The density of the distribution is indicated by the width of each patch, with wider regions, usually in the middle, indicating a higher concentration of data points near those values.

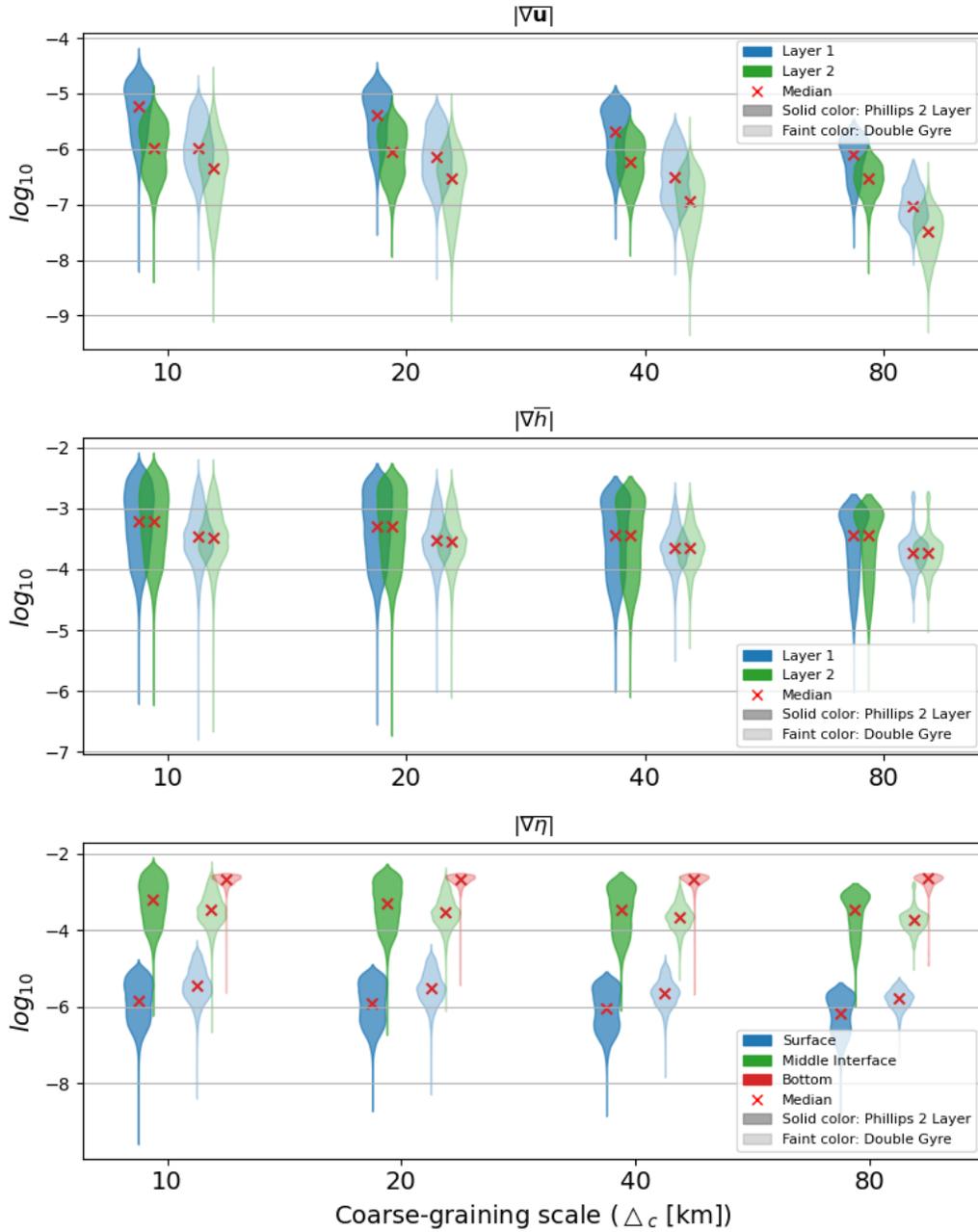


Figure 3. Distribution of input data: Distributions showing the logarithmic magnitude range of different filtered fields, which may be used as input variables for neural network design from both simulations and at different layers and interfaces. Similar to Figure 2 the width of the patch corresponds to values with a higher probability of occurrence.

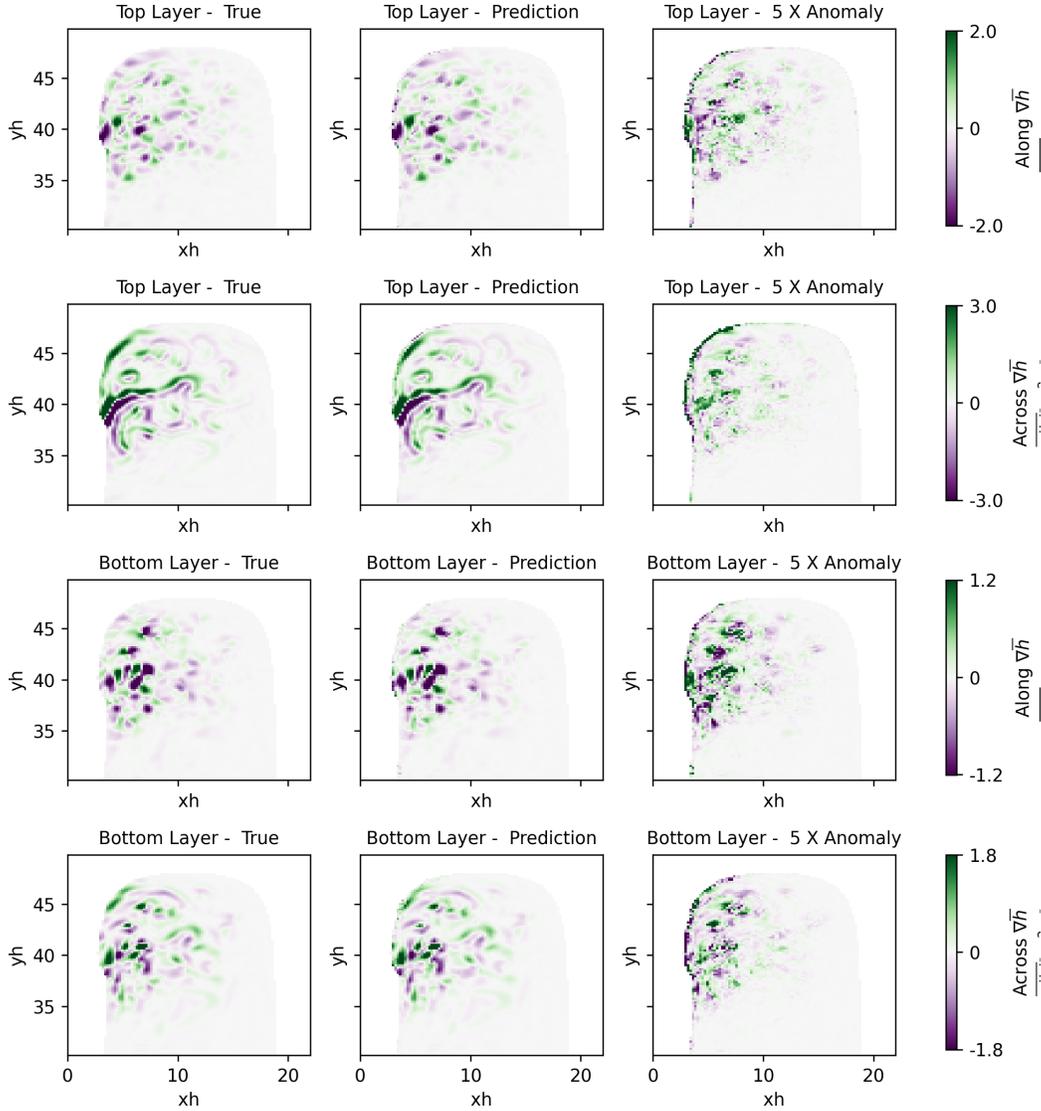


Figure 4. The true (1st column), predicted (2nd column), and prediction error (3rd column) in the along (1st and 3rd row) and across (2nd and 4th row) thickness gradient fluxes for the top (1st and 2nd row) and bottom (3rd and 4th row) layers of the Double Gyre simulations. Here the results are shown for the ML model of the following configuration: 3X3 stencil, trained on data from DG+P2L, and 5090 learnable parameters. These offline skill results are shown for the filter scale of 100km. Note that the prediction error has been multiplied by a factor of 5 to be easily visible on the same scale as the true and predicted fluxes.

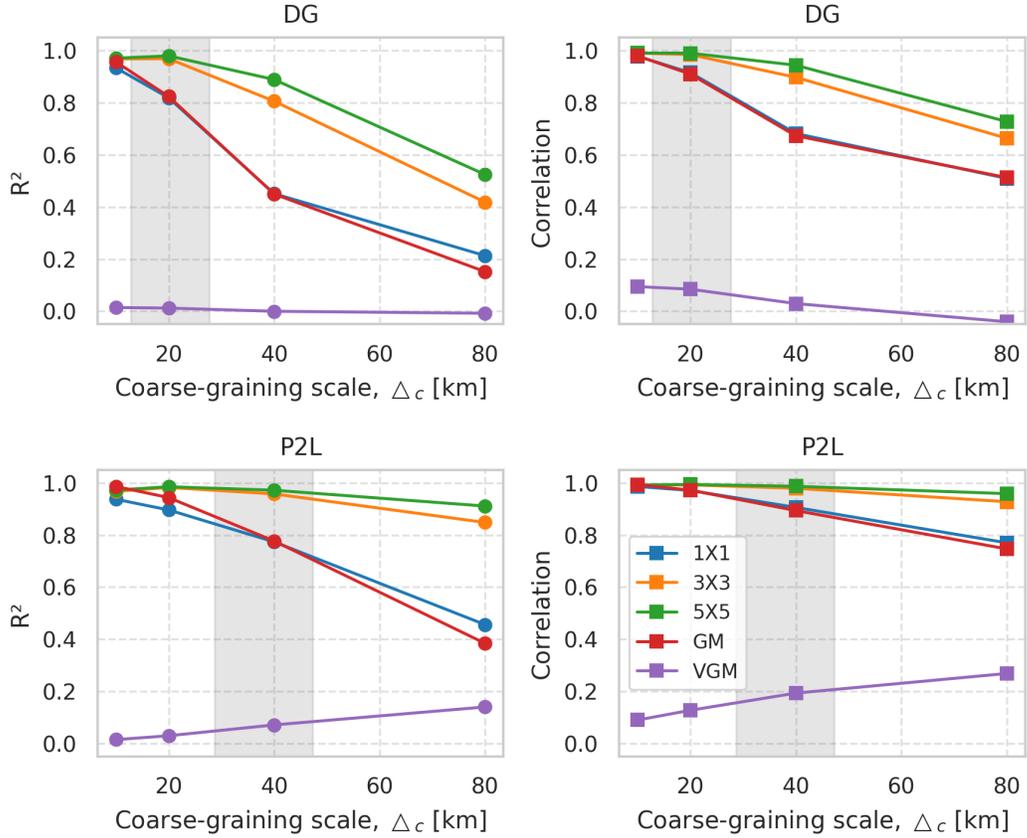


Figure 5. Pointwise offline skill in terms of the R^2 skillscore (left column) and correlation (right column) for 5 different models (indicated in the legend) at different scales (x-axis) and in the double gyre (top) and Phillips 2 layer (bottom) simulations. The metrics were evaluated over both layers and across the full test dataset spanning X years. The gray shaded area indicates the deformation radius range (10 to 90th percentile) in the simulation.

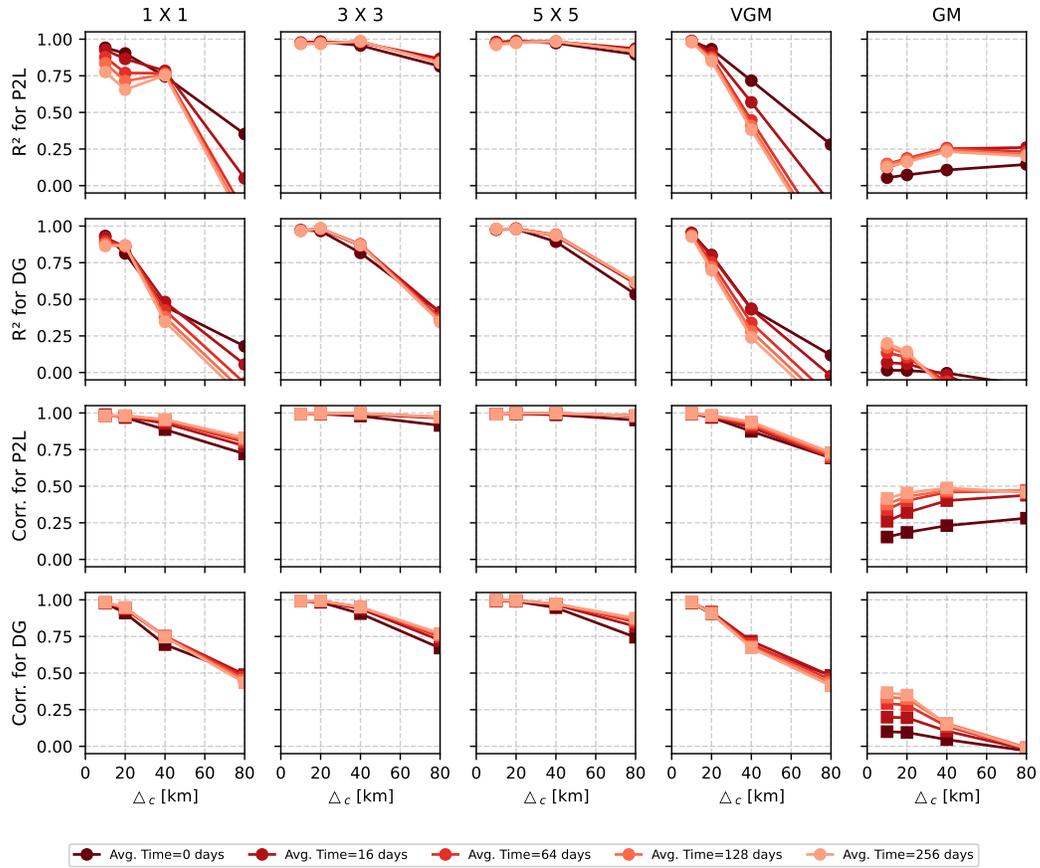


Figure 6. Offline skill in predicting the time average of the subgrid fluxes, using five different models (columns). For all models the skill in predicting the time averaged flux is quantified as the average skill over the two layers and for both the along and across thickness gradient direction.

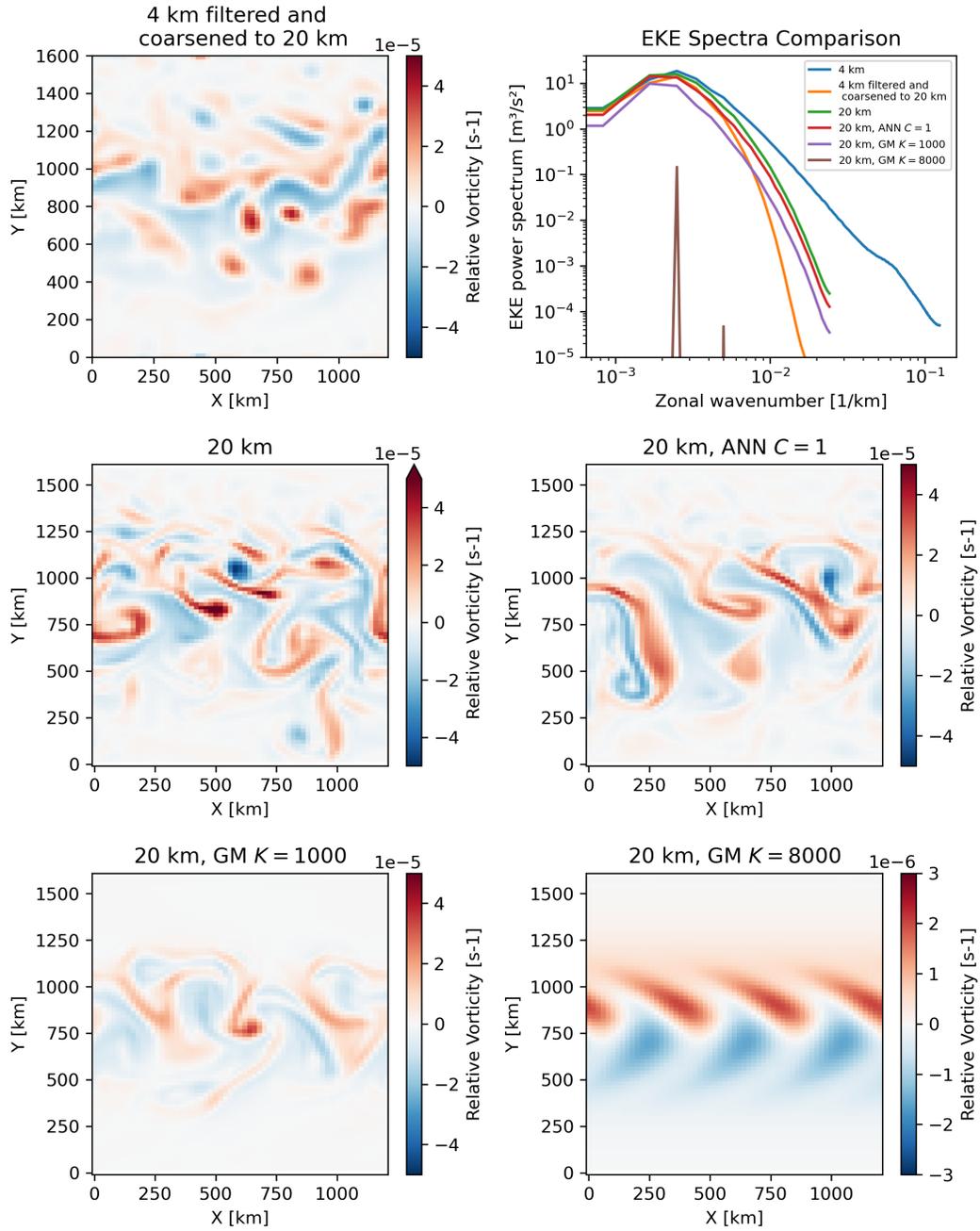


Figure 7. Snapshots of relative vorticity and the EKE spectra from different Phillips 2 layer simulations discussed in section 4.2. Note that the colorbar on the bottom right panel has been adjusted to show the range of values in that simulation.

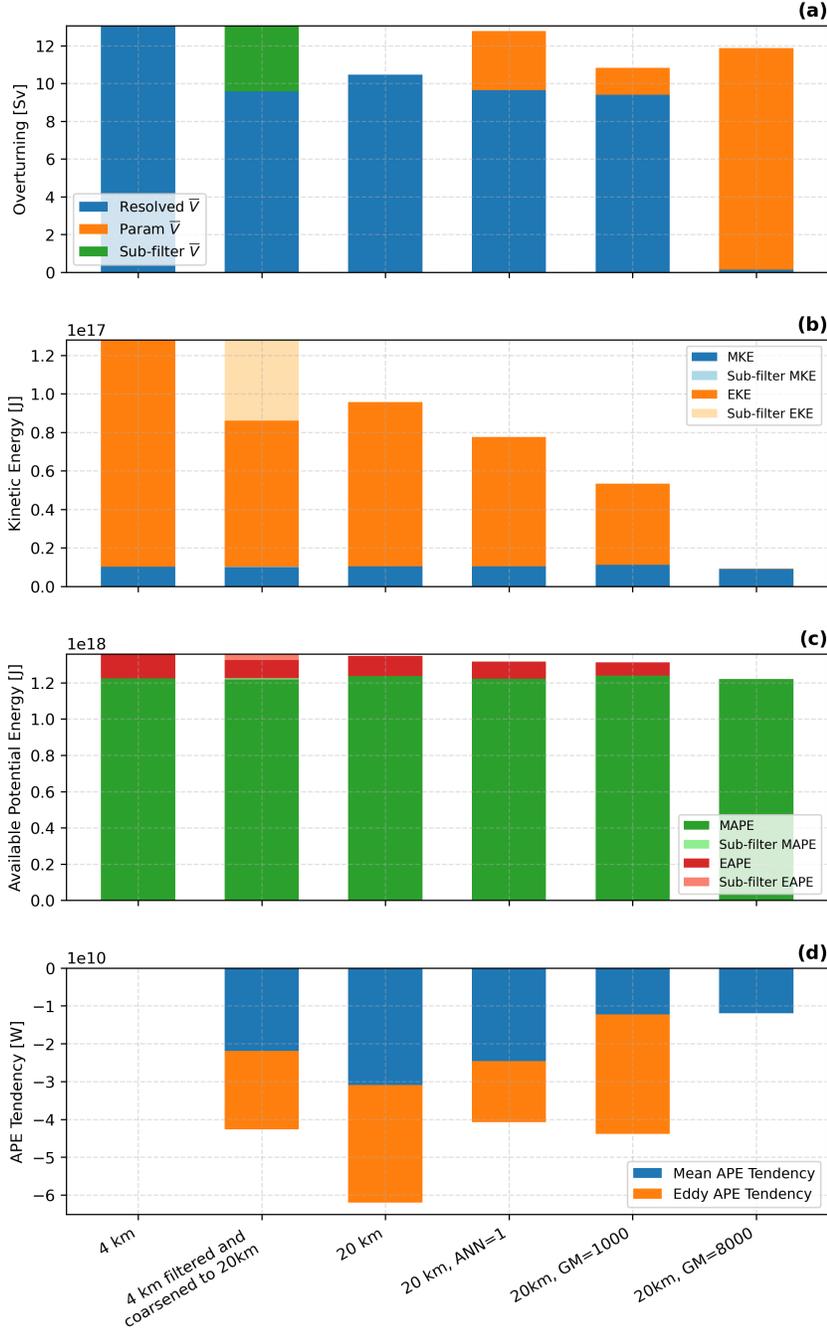


Figure 8. Bulk metrics to evaluate the online skill of the parameterizations in different experiments of Phillips 2 layer simulation (indicated along the x-axis). (a) Peak value of the total meridional transport in the upper layer for the resolved, parameterized and sub-filter components. (b) Kinetic energy and (c) available potential energy of the mean and eddy flow coming from the resolved and SGS contributions. (d) The impact of the parameterized or SGS fluxes on the mean and eddy APE tendency; no bar plot is shown for the high-resolution simulation as in this instance there are no sub-grid or parameterized thickness flux, and the bar for the 20km simulation is calculated using SGS fluxes predicted by a MLP that was not coupled with the resolved fields of the simulation.

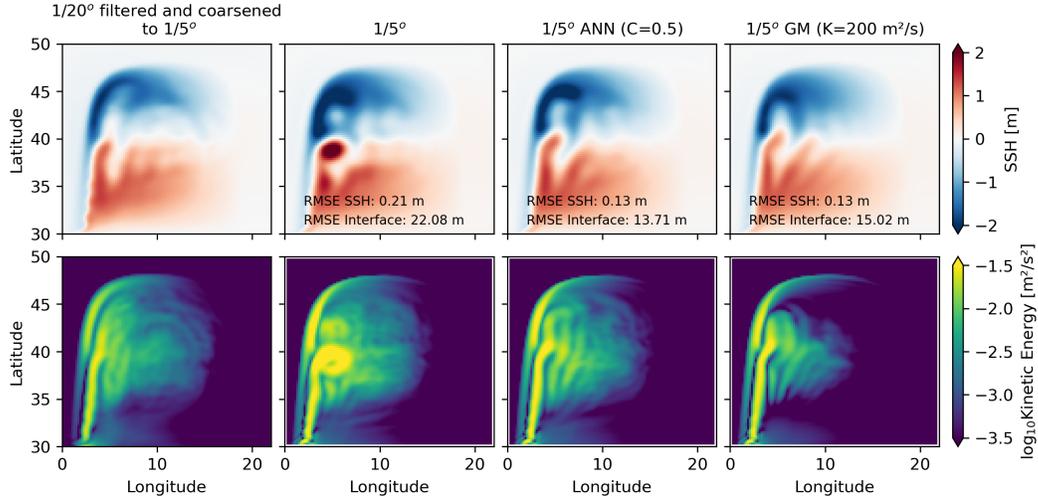


Figure 9. Mean SSH (top row) and EKE (bottom row) the for the Double gyre simulations discussed in section 4.3. The RMSE in the SSH and middle interface height are indicated for the three $1/5^\circ$ resolution simulations.

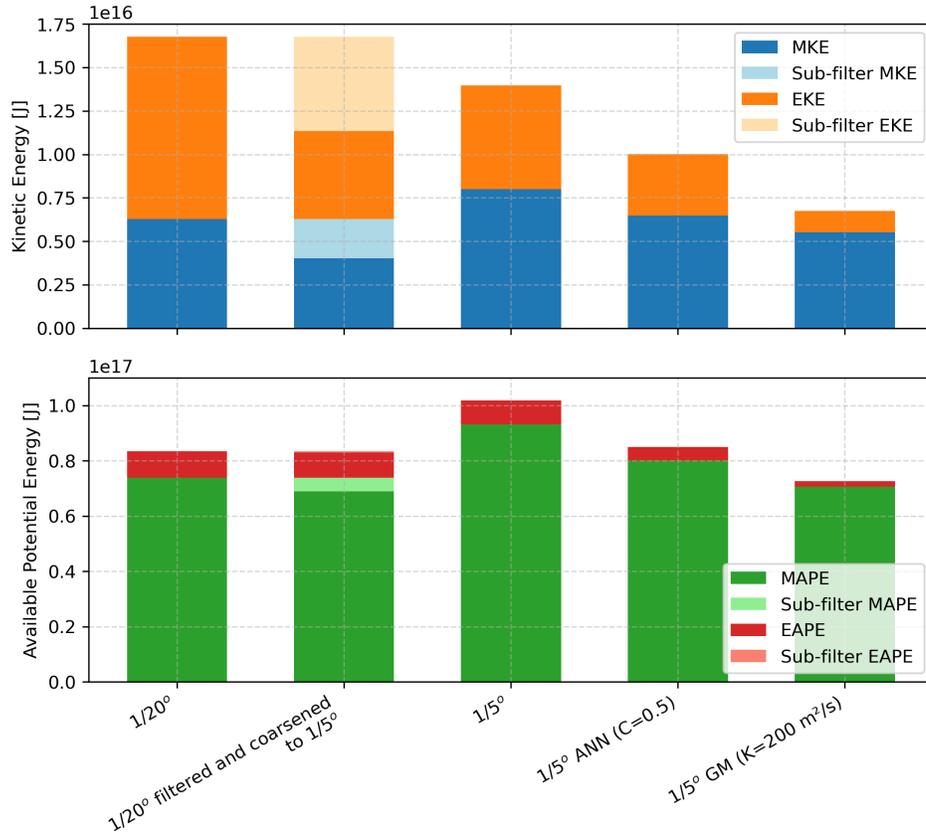


Figure 10. The volume integrated KE (top) and APE (bottom) for the Double Gyre simulations (indicated along x-axis) discussed in section 4.3. The mean, eddy, and sub-filter contributions are indicated.

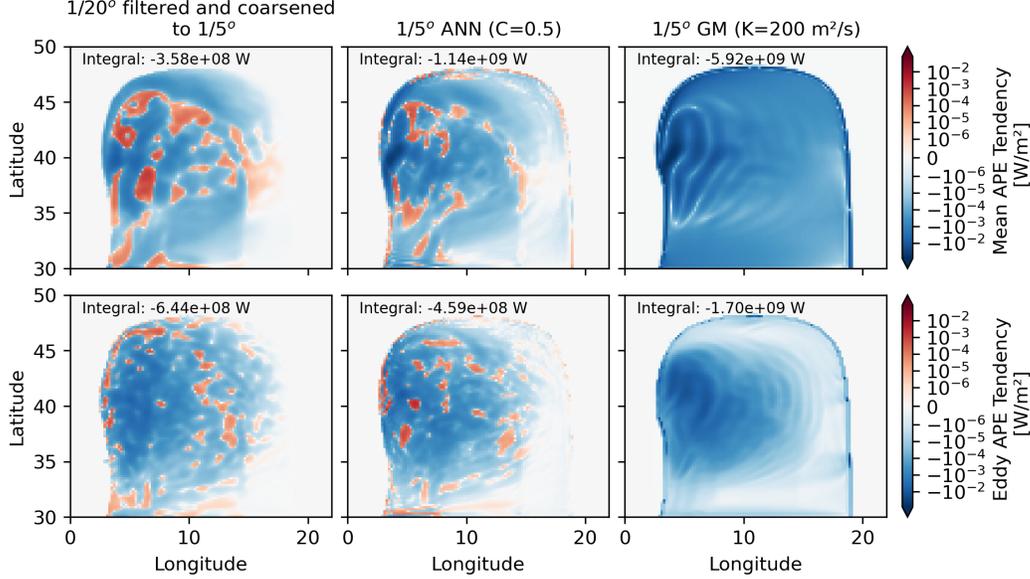


Figure 11. APE tendency exerted on the mean (top) and eddy (bottom) APE by the sub-grid or parameterized fluxes in the Double Gyre simulations discussed in section 4.3. The volume integrated value for each panel is indicated at the top left.

751

Appendix B Eddy driven stream function

The SGS thickness flux corresponds to a volume flux in a layer (\mathbf{F}_n), and the net SGS volume flux below a certain interface can be represented as a stream function (McDougall & McIntosh, 2001):

$$\Psi_{n-1/2} = \sum_{i=N}^n \mathbf{F}_i. \quad (\text{B1})$$

Inversely, the SGS thickness flux in any layer can be expressed in terms of these SGS stream function as:

$$\mathbf{F}_n = \overline{\mathbf{u}_n h_n} - \overline{\mathbf{u}_n h_n} \quad (\text{B2})$$

$$= \Psi_{n-1/2} - \Psi_{n+1/2} \quad (\text{B3})$$

$$= \delta_n \Psi \quad (\text{B4})$$

$$(\text{B5})$$

752

This streamfunction represents a 2D divergent bolus velocity, $\mathbf{u}_n^* = \delta_n \Psi / \overline{h_n}$.

753

754

755

756

757

758

While we have presented the problem entirely in terms of thickness fluxes to be used in layered models (section 2), a comparable representation of SGS buoyancy fluxes in terms of a stream function can be done for depth-level models. When representing the SGS buoyancy fluxes in depth-level models, a 3D non-divergent velocity field is constructed using this streamfunction - the quasi-stokes velocity - with no-flow boundary conditions at the top and bottom may be used.

759 **Appendix C Interface Diffusion and Gent-McWilliams (GM) Param-**
 760 **eterization**

This SGS thickness flux in MOM6 is parameterized by the interface diffusion parameterization, which prescribes the eddy driven streamfunction as,

$$\Psi_{n-1/2}^{GM} = -\kappa^{GM} \nabla \eta_{n-1/2}, \quad 2 \leq n \leq N. \quad (C1)$$

761 In this scheme, the stream function at the top and bottom of the water column are pre-
 762 scribed to be zero ($\Psi_{1/2}^{GM} = \Psi_{N+1/2}^{GM} = 0$), which ensures that the sub-grid thickness
 763 fluxes do not result in any barotropic SGS volume transport. This parameterization al-
 764 ways flattens isopycnal surfaces and reduces APE.

765 This parameterization is a close cousin of the Gent-McWilliams parameterization
 766 (Gent & McWilliams, 1990; Gent et al., 1995), which is designed for use in z-level mod-
 767 els, and also reduces APE of the resolved state.

768 **Appendix D Velocity Gradient Model (VGM)**

769 The velocity gradient model is a structural model, where the goal is to accurately
 770 predict the patterns in the SGS forcing. It is derived based on a Taylor series expansion,
 771 and by keeping only the first term of the expansion (see derivation in Aluie et al. (2022)
 772 or appendix B of Khani and Dawson (2023)).

For thickness fluxes, the VGM predicted flux is:

$$\mathbf{F}_n^{VGM} = C_{VGM} \Delta_c^2 \nabla \bar{\mathbf{u}}_n \nabla \bar{h}_n \quad (D1)$$

where

$$\nabla \bar{\mathbf{u}}_n = \begin{bmatrix} \partial_x \bar{u}_n & \partial_y \bar{u}_n \\ \partial_x \bar{v}_n & \partial_y \bar{v}_n \end{bmatrix}, \quad (D2)$$

$$\nabla \bar{h}_n = \begin{bmatrix} \partial_x \bar{h}_n \\ \partial_y \bar{h}_n \end{bmatrix}, \quad (D3)$$

773 Δ_c is a coarse-graining scale, and C_{VGM} is a scaling coefficient corresponding to the fil-
 774 ter scale. Note that in the presence of topography, we defined our filters such that $\bar{\eta}_b =$
 775 η_b . So, the sub-grid thickness flux in the bottom layer is just the sub-grid deformable
 776 thickness flux.

777 Unlike the GM parameterization, the VGM based expression is not guaranteed to
 778 the APE reducing, and the impact depends non-linearly on the resolved flow.

779 **Appendix E Rotation to thickness gradient frame**

In this work we often rotate all directional quantities: SGS flux vectors, the veloc-
 ity gradient tensor and all thickness gradients, to a thickness gradient frame. This thick-
 ness gradient frame is defined using a set of orthogonal vectors in the horizontal plane.
 The first of these vectors,

$$\hat{\mathbf{T}} = \frac{\nabla \bar{h}_n}{|\nabla \bar{h}_n|} = \frac{\partial_x \bar{h}_n}{|\nabla \bar{h}_n|} \hat{\mathbf{i}} + \frac{\partial_y \bar{h}_n}{|\nabla \bar{h}_n|} \hat{\mathbf{j}}, \quad (E1)$$

points down the thickness gradient. The orthogonal (second) vector, following the right
 hand rule, is defined as,

$$\hat{\mathbf{N}} = \hat{\mathbf{k}} \times \frac{\nabla \bar{h}_n}{|\nabla \bar{h}_n|} = \frac{-\partial_y \bar{h}_n}{|\nabla \bar{h}_n|} \hat{\mathbf{i}} + \frac{\partial_x \bar{h}_n}{|\nabla \bar{h}_n|} \hat{\mathbf{j}}. \quad (E2)$$

The corresponding rotation matrix is defined as

$$\mathbf{R}_n = [\hat{\mathbf{T}} \quad \hat{\mathbf{N}}] = \frac{1}{|\nabla \bar{h}_n|} \begin{bmatrix} \partial_x \bar{h}_n & -\partial_y \bar{h}_n \\ \partial_y \bar{h}_n & \partial_x \bar{h}_n \end{bmatrix}. \quad (\text{E3})$$

This matrix can be used to rotate vector components and tensor components into the thickness gradient frame. Example, the SGS flux vector components in rotated frame can be diagnosed as

$$\widetilde{\mathbf{F}}_n = \mathbf{R}_n^T \mathbf{F}_n, \quad (\text{E4})$$

and the velocity gradient tensor components in can be rotated as,

$$\widetilde{\nabla \mathbf{u}}_n = \mathbf{R}_n^T (\nabla \mathbf{u}_n) \mathbf{R}_n. \quad (\text{E5})$$

When the operation in Equation E4 is applied to the thickness gradient itself, we get the expected result

$$\widetilde{\nabla \bar{h}}_n = |\nabla \bar{h}_n| \hat{\mathbf{T}} + 0 \hat{\mathbf{N}}, \quad (\text{E6})$$

780 since this vector does not have any projection orthogonal to the thickness gradient di-
781 rection by definition.

782 Appendix F ML Model Design Sensitivity

783 There is a vast range of design choices that need to be made when working with
784 machine learning models and designing parameterizations. For example there can be sen-
785 sitivity and interdependence on (i) ML model size and architecture (e.g. we chose MLP
786 here), (ii) training data size, (iii) training data source (which idealized simulation is used
787 to train), (iv) learning rate and optimizers, (v) random parameter initialization, (vi) train-
788 ing targets, (vii) norms optimized in loss functions, (viii) input stencil/ domain of in-
789 fluence, (ix) selection of input features etc. It is infeasible to do a complete search over
790 the entire parameter space and some human intuition is often used to guide the design,
791 here we show the impact of some choices that helped guide our decisions.

792 F1 Impact of different network sizes

793 We want the most skillful predictions at the lowest cost (least number of opera-
794 tions per evaluation of the MLP). We expect that the skill will increase with increasing
795 the number of trainable parameters, but likely saturate beyond a certain point as there
796 may be no more predictive power in the input features left to be extracted. In Figure
797 F1 we quantify this behavior for one particular class of models (these were trained with
798 following choices: trained using data from the double gyre experiment, where all the data
799 was rotated to the thickness gradient frame. Non-dimensional velocity and thickness gra-
800 dients were used as inputs. Mean absolute error was used for the loss, with dimensional
801 fluxes (equation 4) as output targets. Both inputs and outputs were normalized using
802 order of magnitude estimates. We used the model snapshots 0-640 for training, 672-736
803 for evaluation and 736-800 for testing. Adam optimizer was used with a learning rate
804 of 0.01, and training was continued till the relative improvement in the loss saturated
805 within a relative tolerance of 0.01 (1% error) for at least 10 epochs).

806 The model skill, measured as the R2 value, improves as the number of parameters
807 increase, and this is particularly apparent at larger filter scales and when the model sten-
808 cil is wider. However, since skill is already very high at filter scales of 50 and 100 km this
809 effect is minor, and at larger filter scales the skill seems to asymptote to its maximum
810 value approximately around 10K parameters.

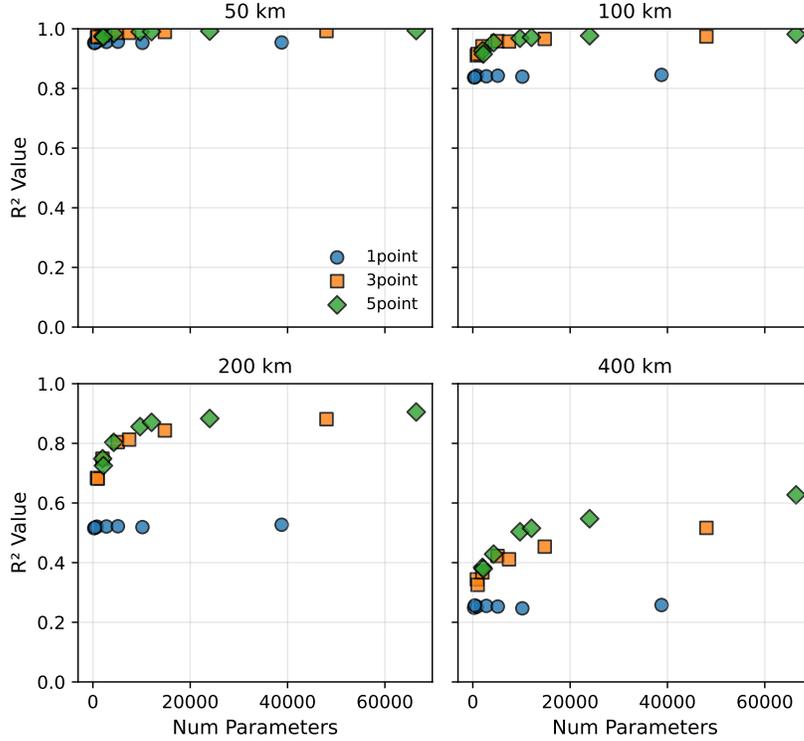


Figure F1. ML model skill, defined as the R2 value, as a function of number of parameters for different filter scales. The details ML model being tested is described in F1.

811

F2 Impact of training data size

812

To test the impact that the amount of training data has, we experimented with the same setup as the one used in the previous section with varying amount of data. We fixed all aspects, and performed the test in the case where the stencil size is 3X3 and ML model has 3386 parameters (model shape 54,36,36,2). We created batches of 16 model snapshots, and tested the effect that increasing the number of batches had. As shown in Figure F2 the model skill increases upto about 128 batches (2048 model snapshots), but saturates past that point. This led us to use this data volume for training the model used in the main text.

819

820

F3 Impact of different loss functions and targets

821

Non-dimensionalization of the sub-grid fluxes leads to a significant collapse in the distributions (see Figure 2), but also produces outliers (due to possibility of division by small velocity or thickness gradients). We have two choices of loss during training, either training on the dimensional fluxes $\mathcal{L}^{dim} = \|\mathbf{F} - \mathbf{F}^{pred}\| = \|\mathbf{F} - (\Delta^2 |\nabla \mathbf{u}| |\nabla h|) f_{\theta}(\cdot)\|$ or training on non-dimensional fluxes $\mathcal{L}^{non-dim} = \|\mathbf{F} / (\Delta^2 |\nabla \mathbf{u}| |\nabla h|) - f_{\theta}(\cdot)\|$. Both these forms should give us the same functional representation if a unique function $f_{\theta}(\cdot)$ exists, but in the more realistic situation where $f_{\theta}(\cdot)$ is an approximation and we want it to equally balance data coming from many regimes (different simulations, varying energy levels at different depths, etc) it would seem that $\mathcal{L}^{non-dim}$ might be a better choice. Along with this we also compared the use of mean square error (MSE) vs mean absolute error (MAE), since MAE is generally considered to be more tolerant to outliers.

831

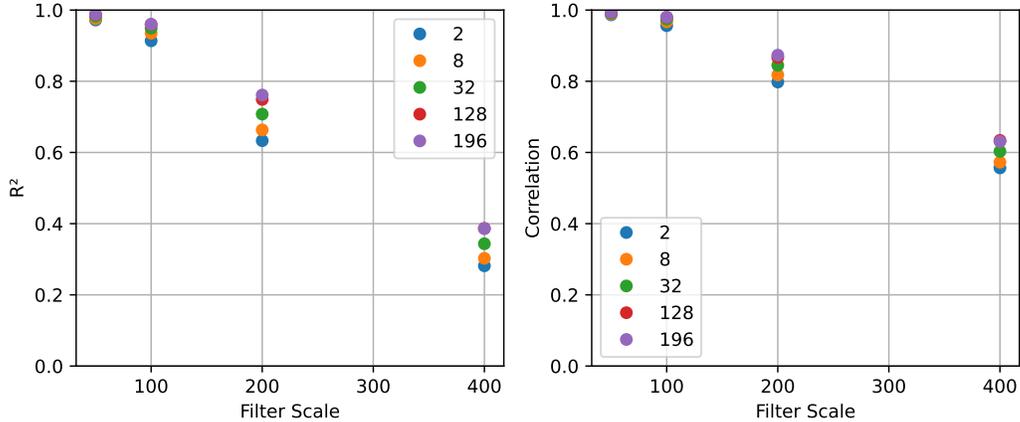


Figure F2. ML model skill, defined as the R² value, as a function of filter scales for different number of data batches used for training.

832 As seen in Figure F3, the impact of these choices is relatively minor in the model
 833 considered in the previous section. Generally, the models trained using MAE seem to
 834 go better, except at the largest filter scales - where the model trained using non-dimensionalized
 835 outputs and MSE does better. However, this model is the worst performer at all other
 836 scales. The models trained using MAE also generally perform well across a wide range
 837 of scales (bottom panel). To our surprise, the impact of this choice was smaller than we
 838 had expected. For the main study we used $\mathcal{L}^{non-dim}$ along with MAE.

839 **F4 Impact of choice of training data**

840 We want to build ML models that are easily generalizable and training data ag-
 841 nostic. For example, a model trained on data from the P2L simulation should be able
 842 to make good predictions on data from DG experiment, and vice versa. Developing ap-
 843 propriate non-dimensionalizations, as highlighted in section 3, is a step in this direction.
 844 In Figure F4 we show that the model trained using data from both simulations simul-
 845 taneously generally performs the best or close to the best, which is why we chose this
 846 training strategy for the model in the main text.

847 It is worth noting that in fact even models trained on a single experiment, perform
 848 relatively well when tested on the other experiment. This is a result of the design choices
 849 we made. During the initial phase of our development, when we trained models to have
 850 the form of equation 4, the skill on unseen data was extremely poor (not shown), also
 851 see Perezhugin et al. (submitted).

852 **F5 Impact of other aspects**

853 The random seed, which sets **the random initializations** of the ML model weights,
 854 seems to have a small impact on the skill. The model skill averaged across all filter scales,
 855 for the case discussed in the above section, varied between values of 0.7 - 0.725 over dif-
 856 ferent random seeds. Since this effect seems small, in the main text we use only a sin-
 857 gle trained model.

858 The **learning rate**, similarly has a minor impact on the final skill but impacts the
 859 nature of decrease in the loss. Large learning rates (~ 0.1) lead to noisy training, while
 860 smaller rates (< 0.005) take too many epochs to train. We found that a learning rate

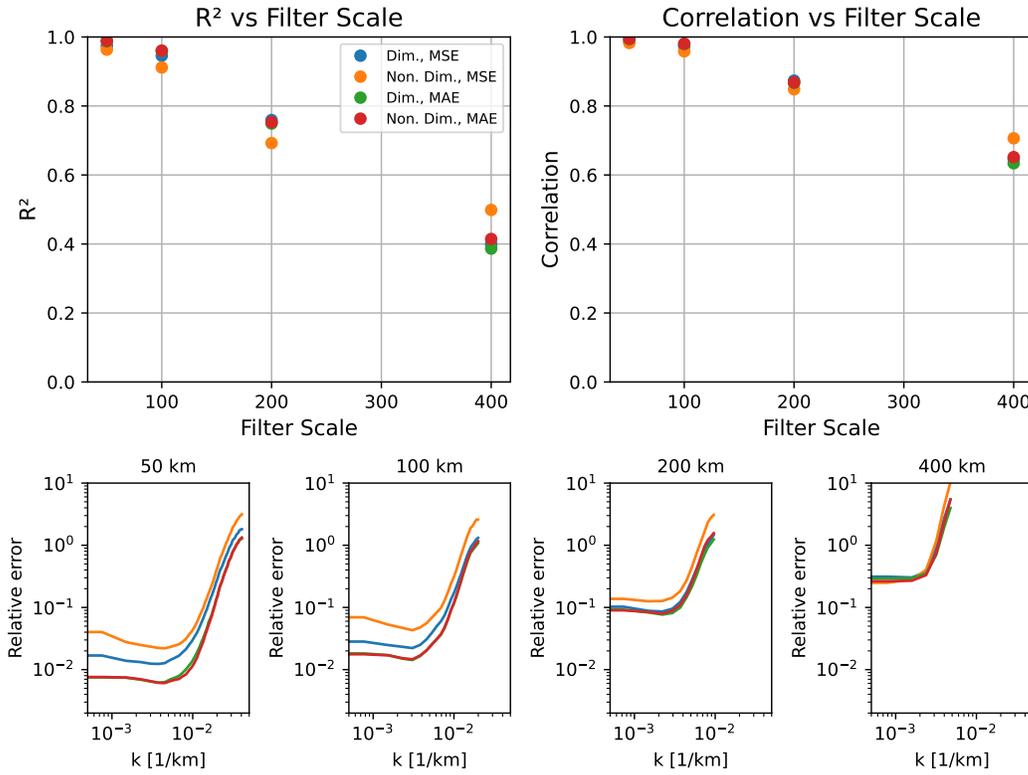


Figure F3. Top row shows the model skill in terms of R² and correlation as a function of filter scale for different choices of loss function. Bottom two rows show the relative error, which is defined as the power spectrum of the (truth - prediction) divided by the power spectrum of the truth; value of greater than 1 implies that the variance in the anomaly field is greater than in the truth.

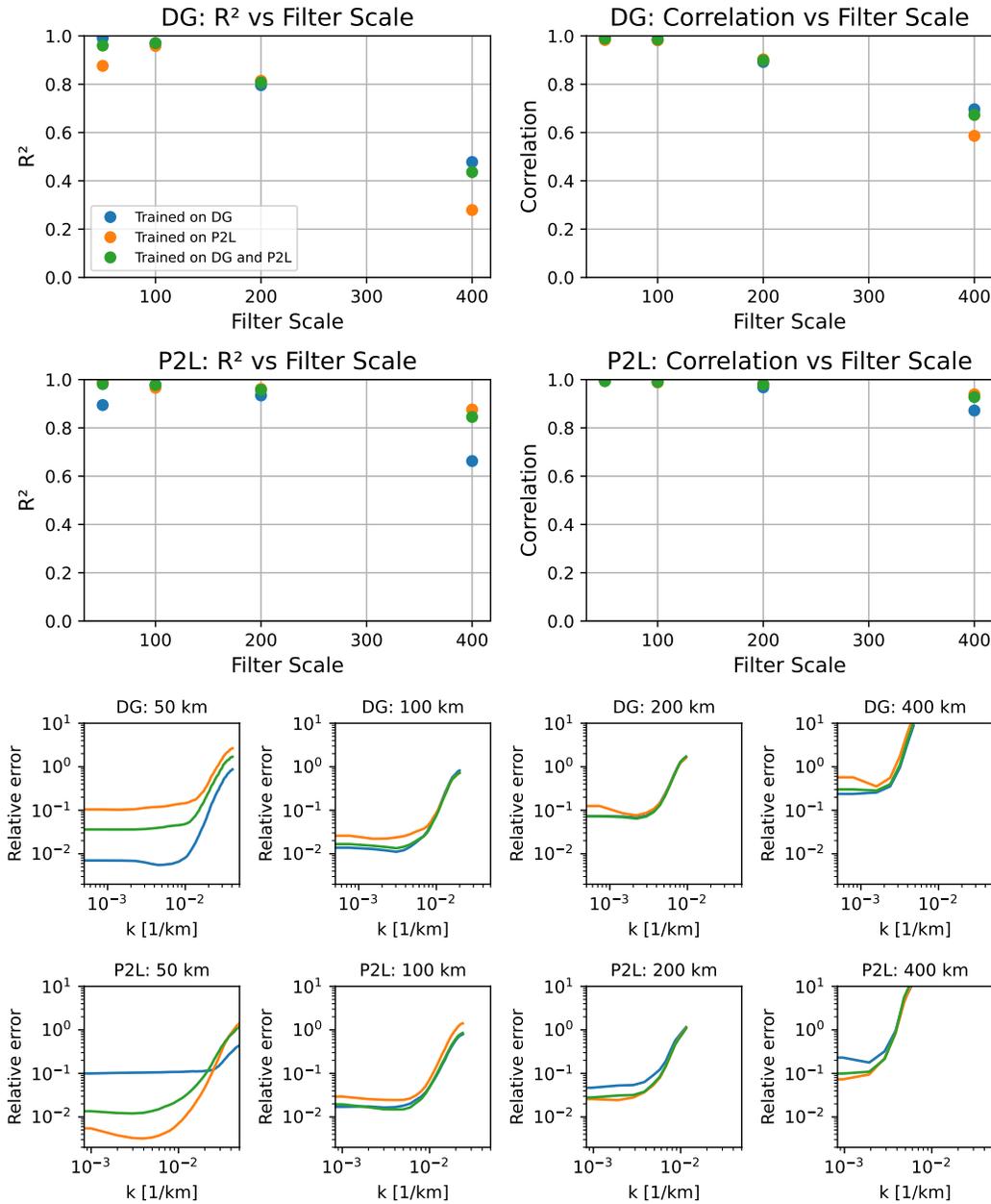


Figure F4. Top two rows show the model skill in terms of R² and correlation as a function of filter scale for different choices of training datasets. Bottom two rows show the relative error, which is defined as the power spectrum of the (truth - prediction) divided by the power spectrum of the truth; value of greater than 1 implies that the variance in the anomaly field is greater than in the truth. The title of the panels indicate the dataset on which testing is done, and the legend indicated the dataset that is used for training.

861 around 0.01 was optimal in not taking too many epochs and not being too noisy, which
 862 is what we used for the models discussed in the main text.

863 Appendix G Simulation evaluation metrics

864 Here we describe the metrics used to assess the physical properties of our simula-
 865 tions.

866 G1 Overturning Circulation

The overturning circulation ($\bar{V}_n(y)$) is defined as the volume transport across a lon-
 gitude band in a particular layer,

$$\bar{V}_n(y) = \oint \overline{v_n h_n^{-t}} dx + \oint \overline{F_n^{yt}} dx, \quad (\text{G1})$$

867 where $F_n^y(x, y, z, t)$ is the sub-grid or sub-filter meridional flux in layer n , $\oint(\cdot)dx$ corre-
 868 sponds to a zonal integral over the full longitudinal domain, and $\overline{(\cdot)^t}$ indicates a time av-
 869 erage. The first term on the RHS corresponds to the resolved overturning and includes
 870 the contribution from both the mean and the variable parts of the flow, and the second
 871 term corresponds to the parameterized overturning.

872 G2 Eddy Kinetic Energy Spectrum

The EKE spectrum provides an effective measure of the flow variability at differ-
 ent scales. Here, we define it as the zonal wavenumber (k) power spectrum of the $\mathbf{u}'_n(x, y, t) =$
 $\mathbf{u}_n(x, y, t) - \bar{\mathbf{u}}_n^t(x, y)$,

$$EKE_n(k) = \frac{1}{2} (\overline{|u'_n(k, y, t)|^2}^{t,y} + \overline{|v'_n(k, y, t)|^2}^{t,y}), \quad (\text{G2})$$

where $\overline{(\cdot)^{t,y}}$ is a time and meridional average and $u'_n(k, y, t)$ is the zonal Fourier trans-
 form of $u'_n(x, y, t)$. By Parseval's theorem we have

$$\sum_k EKE_n(k) = \frac{1}{2} \int (\overline{|u'_n(x, y, t)|^2}^{t,y} + \overline{|v'_n(x, y, t)|^2}^{t,y}) dx, \quad (\text{G3})$$

873 which shows the scale-wise decomposition aspect of the zonal wavenumber EKE spec-
 874 trum in the x-direction. We used the xrft package (<https://xrft.readthedocs.io/>)
 875 for this analysis.

876 G3 Integral Kinetic and Available Potential Energies

The volume integrated KE in Joules is defined as:

$$KE = \frac{1}{2} \sum_k \int \rho_0 |\mathbf{u}_k|^2 h_k dx dy \quad (\text{G4})$$

877 The KE of the mean flow ($\bar{\mathbf{u}}_k^t, \bar{h}_k^t$) is referred to as MKE , and the EKE is defined
 878 as $EKE = \overline{KE}^t - MKE$.

The volume integrated APE in Joules is defined as:

$$APE = \frac{1}{2} \sum_k \int \rho_0 g'_{k-1/2} (\eta_{k-1/2} - \eta_{k-1/2}^{ref})^2 dx dy \quad (\text{G5})$$

879 The $MAPE$ is the APE of the mean state ($\overline{\eta_{k-1/2}^t}$), and the $EAPE$ is given by
 880 $EAPE = \overline{APE}^t - MAPE$

881 For filtered and coarsened flow ($\overline{\mathbf{u}_k^\Delta}, \overline{h_k^\Delta}$) we can also define the total kinetic energy (KE^Δ), which has its corresponding mean (MKE^Δ , for $\overline{\mathbf{u}_k^{\Delta,t}}, \overline{h_k^{\Delta,t}}$) and eddy (EKE^Δ) components. Accordingly, we also define the SGS kinetic energies for the total ($KE - MKE^\Delta$), mean ($MKE - MKE^\Delta$), and eddy ($EKE - EKE^\Delta$) flow. Similarly APE and its components can be defined for the filtered and coarsened flow and the SGS contribution.

887 **G4 Available potential energy tendency due to sub-filter or parameterized fluxes**
888

Following (Loose, Marques, et al., 2023), we can compute the impact that the SGS fluxes have on volume integrated APE tendency as,

$$(\partial_t APE)^{SF} = \sum_{n=1}^N \rho_0 \mathbf{F}_n \cdot \nabla M_n, \quad (G6)$$

889 where \mathbf{F}_n are the SGS fluxes and M_n is the dynamic pressure of the resolved.

In a 2 layer fluid we have,

$$(\partial_t APE)^{SF} = \rho_0 g_{1/2}^r (\mathbf{F}_1 + \mathbf{F}_2) \cdot \nabla \eta_{1/2} + \rho_0 g_{3/2}^r \mathbf{F}_2 \cdot \nabla \eta_{3/2}. \quad (G7)$$

890 Notably, only the last term on the RHS, which arises due to SGS fluxes in the bottom layer, makes a significant contribution to the APE tendency, since $|\nabla \eta_{3/2}| \gg |\nabla \eta_{1/2}|$.
891 The barotropic contribution to APE, first term on RHS, is small.
892

We decompose this tendency into the time mean and eddy contribution, by defining the contribution from the mean as

$$(\partial_t APE)^{SF,mean} = \rho_0 g_{3/2}^r (\overline{\mathbf{F}_2}^t \cdot \nabla \overline{\eta_{3/2}}^t), \quad (G8)$$

893 where the barotropic contribution is neglected. The eddy contribution is defined as $(\partial_t APE)^{SF,eddy} =$
894 $\frac{(\partial_t APE)^{SF,t} - (\partial_t APE)^{SF,mean}}{(\partial_t APE)^{SF,t} - (\partial_t APE)^{SF,mean}}$.

895 **Appendix H Online sensitivity to grid-size and parameterization coefficients**
896

897 Simulation of oceanic mesoscale turbulence is relatively sensitive to how well the first baroclinic deformation radius (shown in Figure H1 for the two simulations considered here) is resolved. In this study we chose our coarsening scale and simulation grid-size to lie in a range that encompasses the deformation radius, so that skill of the new parameterization at both the eddy permitting and non-eddying resolutions can be investigated. Here we provide details of the sensitivity of the coarse simulations to both grid spacing and parameterization coefficients.
898
899
900
901
902
903

904 **H1 Phillips 2 Layer**

905 For P2L simulation we tested the sensitivity of the overturning transport to the grid size and parameterization coefficients, shown in Figure H2.
906

907 In this case, the unparameterized low resolution simulation always has lower overturning transport relative to the HR simulation. The response of the parameterized flux contribution to the GM diffusivity is almost linear and insensitive to grid size, and adjusting the GM diffusivity allows us to increase the total overturning transport to the appropriate value. However, this is only achieved when the diffusivity is relatively large ($\kappa_{GM} \sim 8,000-10,000 \text{ m}^2/\text{s}$), which is also a parameter regime in which the resolved eddy flow has been entirely suppressed. Consequently the contribution of the resolved
908
909
910
911
912
913

914 flow to the overturning has been completely eliminated and the entire contribution comes
 915 from parameterized fluxes (Figure H2 second column). There seems to be no lower dif-
 916 fusivity value at which the appropriate overturning can be achieved with only marginal
 917 impact to the resolved flow.

918 In contrast, the simulations with the ANN based parameterization are usually able
 919 to achieve the appropriate level of overturning with minimal damage to the resolved flow.
 920 The response of the parameterized flux to the grid size and parameterization amplifica-
 921 tion coefficient (C_{ANN}) is non-linear. At grid sizes of 10 and 20 km, which are smaller
 922 than the deformation radius everywhere in the domain, the appropriate overturning is
 923 achieved at $C_{ANN} = 1$ and with minimal damage to the resolved overturning. At a grid
 924 size of 80 km, which is larger than the deformation radius everywhere in the domain, the
 925 appropriate overturning is achieved at $C_{ANN} = 2$. At the intermediate grid size of 40 km,
 926 the ANN parameterization does produce improvements to overall overturning with lit-
 927 tle damage to the resolved flow (also at $C_{ANN} = 1$), but is unable to achieve similar
 928 overturning to the HR simulation. We anticipate that combining the thickness flux pa-
 929 rameterization with a momentum flux parameterization may be able to produce further
 930 improvements at these intermediate resolutions.

931 H2 Double Gyre

932 Unlike the P2L simulation, in the DG case there is no overturning circulation. In
 933 the context of the thickness flux parameterization, the error in the MAPE is the most
 934 relevant, which is also linked to the error in mean SSH. We tested the sensitivity of these
 935 quantities and a few others to the parameterization coefficients and grid-size (Figure H3).

936 For the GM parameterization, a diffusivity of about $200 \text{ m}^2/\text{s}$ is able to achieve
 937 the appropriate MAPE at all resolutions, which similar to the P2L also comes an almost
 938 complete damping of the resolved eddies. The ANN is able to achieve the best MAPE
 939 with $C_{ANN} \sim 0.5-0.75$, and this happens with relatively less deterioration in the re-
 940 solved eddies.

Δ_c [km]	κ_{GM} [m^2/s]		C_{VGM} [nondim]	
	P2L	DG	P2L	DG
10	137	40	0.110	0.075
20	596	105	0.077	0.072
40	2287	192	0.077	0.067
80	6266	108	0.091	0.048

Table H1. Estimated GM diffusivity (κ_{GM}) and VGM coefficient (C_{VGM}) from the data using least squares fitting, for both the Phillips 2-Layer (P2L) and Double Gyre (DG) setups.

941 References

- 942 Abernathey, R., & Wortham, C. (2015). Phase speed cross spectra of eddy heat
 943 fluxes in the eastern pacific. *Journal of Physical Oceanography*, *45*(5), 1285–
 944 1301.
- 945 Adcroft, A., Anderson, W., Balaji, V., Blanton, C., Bushuk, M., Dufour, C. O., ...
 946 others (2019). The gfdl global ocean and sea ice model om4. 0: Model descrip-
 947 tion and simulation features. *Journal of Advances in Modeling Earth Systems*,
 948 *11*(10), 3167–3211.

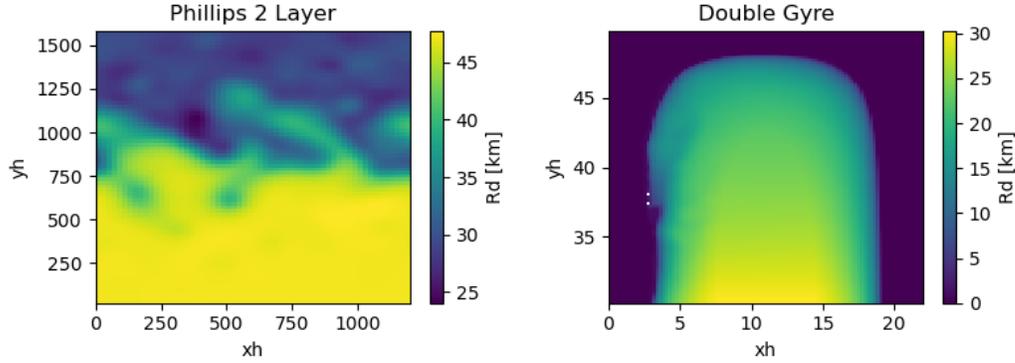


Figure H1. First baroclinic deformation radii for P2L (left) and DG (right). Note the different range of scales for the two colobars.

- 949 Aluie, H., Hecht, M., & Vallis, G. K. (2018). Mapping the energy cascade in the
 950 north atlantic ocean: The coarse-graining approach. *Journal of Physical*
 951 *Oceanography*, *48*(2), 225–244.
- 952 Aluie, H., Rai, S., Yin, H., Lees, A., Zhao, D., Griffies, S. M., . . . Shang, J. K.
 953 (2022). Effective drift velocity from turbulent transport by vorticity. *Physical*
 954 *Review Fluids*, *7*(10), 104601.
- 955 Beucler, T., Pritchard, M., Yuval, J., Gupta, A., Peng, L., Rasp, S., . . . others
 956 (2021). Climate-invariant machine learning. *arXiv preprint arXiv:2112.08440*.
- 957 Bodner, A., Balwada, D., & Zanna, L. (2023). A data-driven approach for parame-
 958 terizing submesoscale vertical buoyancy fluxes in the ocean mixed layer. *arXiv*
 959 *preprint arXiv:2312.06972*.
- 960 Bracco, A., Brajard, J., Dijkstra, H. A., Hassanzadeh, P., Lessig, C., & Monteleoni,
 961 C. (2025). Machine learning for the physics of climate. *Nature Reviews*
 962 *Physics*, *7*(1), 6–20.
- 963 Brenowitz, N. D., & Bretherton, C. S. (2018). Prognostic validation of a neural net-
 964 work unified physics parameterization. *Geophysical Research Letters*, *45*(12),
 965 6289–6298.
- 966 Chelton, D. B., DeSzoeke, R. A., Schlax, M. G., El Naggar, K., & Siwertz, N.
 967 (1998). Geographical variability of the first baroclinic rossby radius of de-
 968 formation. *Journal of Physical Oceanography*, *28*(3), 433–460.
- 969 Eden, C., & Greatbatch, R. J. (2008). Towards a mesoscale eddy closure. *Ocean*
 970 *Modelling*, *20*(3), 223–239.
- 971 Ferrari, R., Griffies, S. M., Nurser, A. G., & Vallis, G. K. (2010). A boundary-value
 972 problem for the parameterized mesoscale eddy transport. *Ocean Modelling*,
 973 *32*(3-4), 143–156.
- 974 Ferrari, R., McWilliams, J. C., Canuto, V. M., & Dubovikov, M. (2008). Parame-
 975 terization of eddy fluxes near oceanic boundaries. *Journal of Climate*, *21*(12),
 976 2770–2789.
- 977 Ferrari, R., & Wunsch, C. (2009). Ocean circulation kinetic energy: Reservoirs,
 978 sources, and sinks. *Annual Review of Fluid Mechanics*, *41*, 253–282.
- 979 Ferreira, D., Marshall, J., & Heimbach, P. (2005). Estimating eddy stresses by fit-
 980 ting dynamics to observations using a residual-mean ocean circulation model
 981 and its adjoint. *Journal of Physical Oceanography*, *35*(10), 1891–1910.
- 982 Gent, P. R. (2011). The gent-mcwilliams parameterization: 20/20 hindsight. *Ocean*
 983 *Modelling*, *39*(1-2), 2–9.
- 984 Gent, P. R., & McWilliams, J. C. (1990). Isopycnal mixing in ocean circulation mod-
 985 els. *Journal of Physical Oceanography*, *20*(1), 150–155.

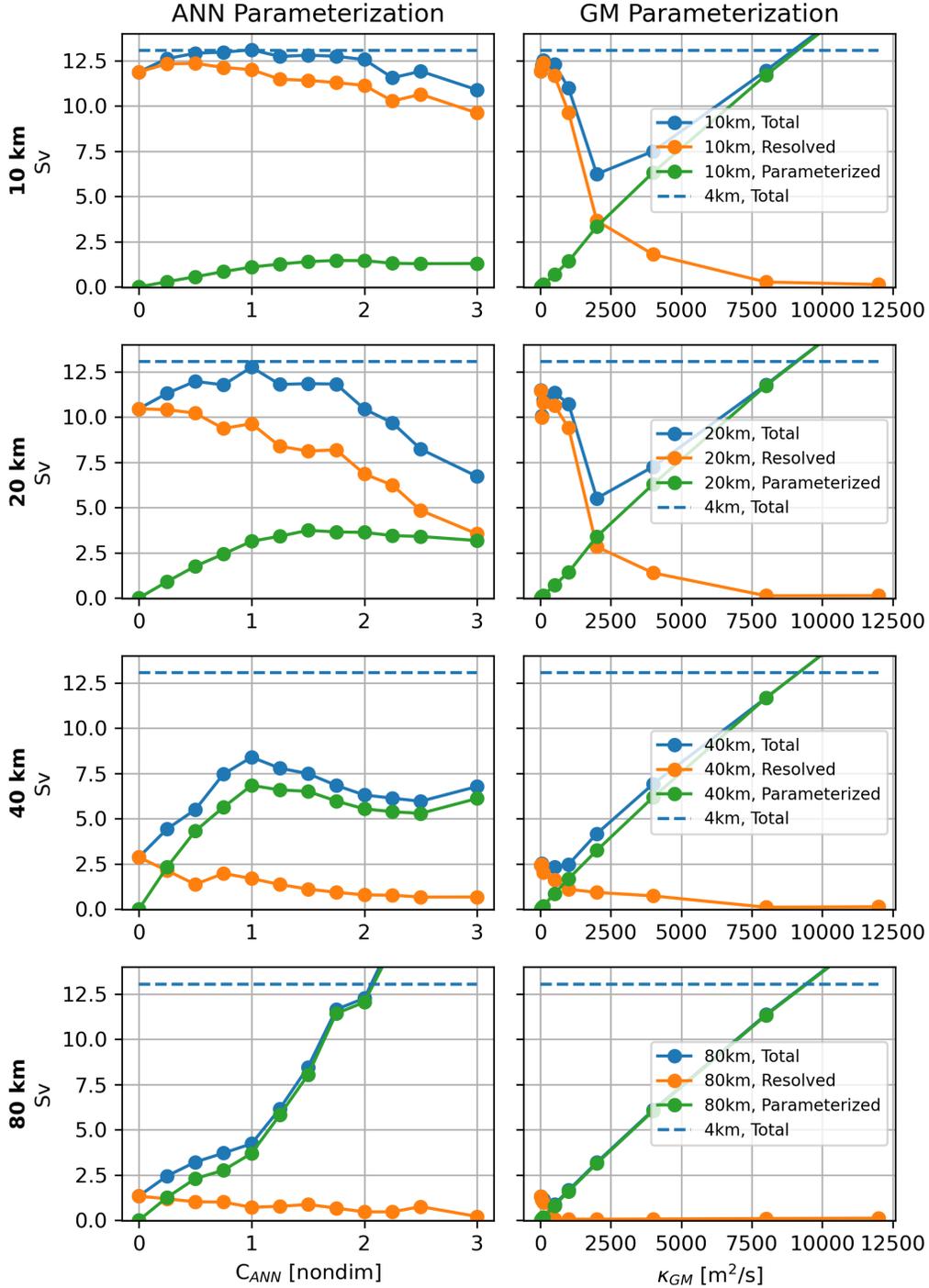


Figure H2. Sensitivity of the overturning circulation and its resolved and parameterized components to the MLP coefficient (left) and the GM diffusivity (right) for the Phillip 2 layer simulations at different grid-sizes.

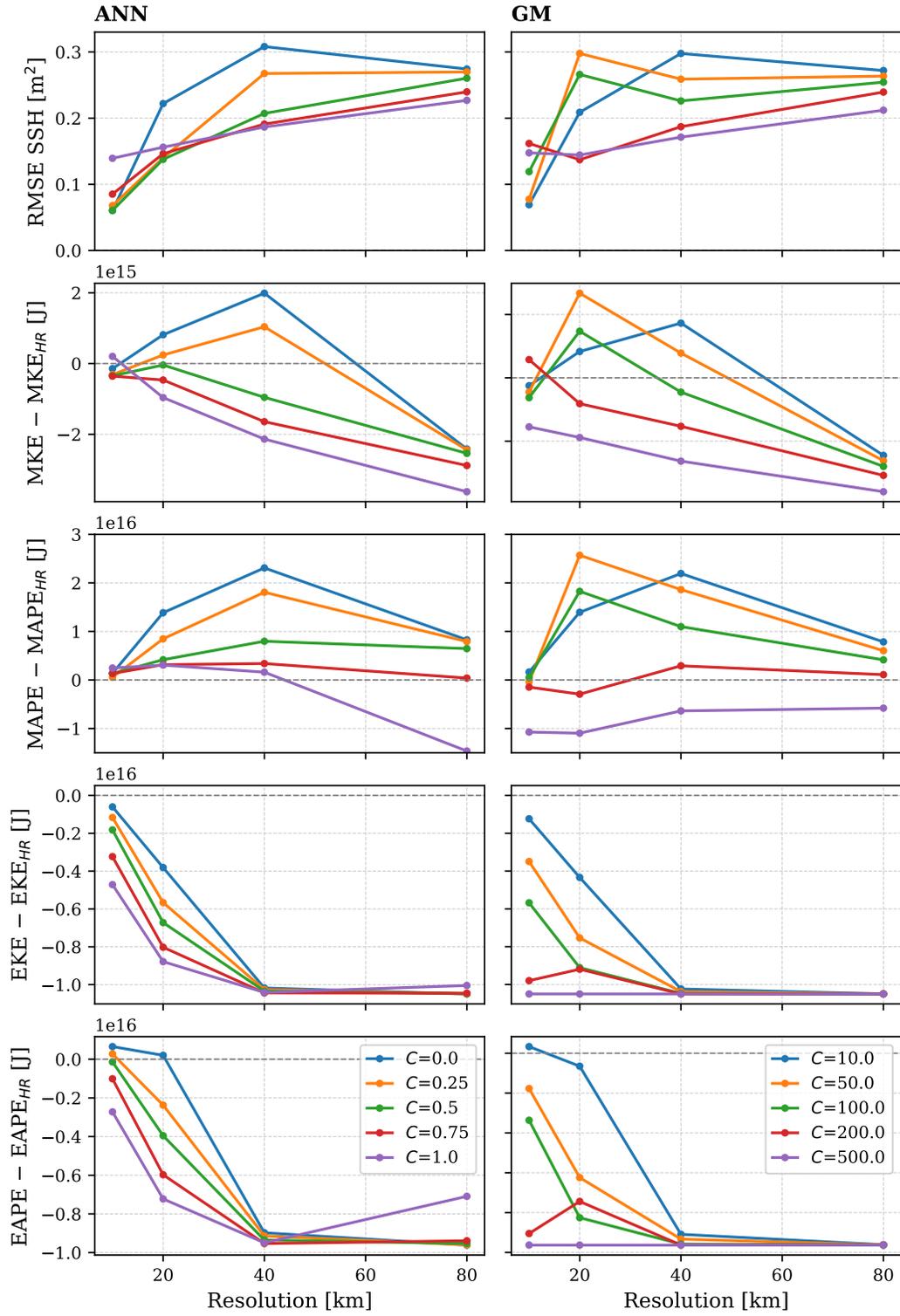


Figure H3. Sensitivity to the MLP coefficient (left) and the GM diffusivity (right) for the Double Gyre simulations at different resolutions.

- 986 Gent, P. R., Willebrand, J., McDougall, T. J., & McWilliams, J. C. (1995). Parame-
 987 terizing eddy-induced tracer transports in ocean circulation models. *Journal of*
 988 *physical oceanography*, 25(4), 463–474.
- 989 Hallberg, R. (2013). Using a resolution function to regulate parameterizations of
 990 oceanic mesoscale eddy effects. *Ocean Modelling*, 72, 92–103.
- 991 Jansen, M. F., Adcroft, A. J., Hallberg, R., & Held, I. M. (2015). Parameterization
 992 of eddy fluxes based on a mesoscale energy budget. *Ocean Modelling*, 92, 28–
 993 41.
- 994 Khani, S., & Dawson, C. N. (2023). A gradient based subgrid-scale parameteriza-
 995 tion for ocean mesoscale eddies. *Journal of Advances in Modeling Earth Sys-*
 996 *tems*, 15(2), e2022MS003356.
- 997 Killworth, P. D. (2001). Boundary conditions on quasi-stokes velocities in parame-
 998 terizations. *Journal of physical oceanography*, 31(4), 1132–1155.
- 999 Lai, C.-Y., Hassanzadeh, P., Sheshadri, A., Sonnewald, M., Ferrari, R., & Balaji, V.
 1000 (2024). Machine learning for climate physics and simulations. *Annual Review*
 1001 *of Condensed Matter Physics*, 16.
- 1002 Loose, N., Abernathey, R., Grooms, I., Busecke, J., Guillaumin, A., Yankovsky, E.,
 1003 ... others (2022). Gcm-filters: A python package for diffusion-based spatial
 1004 filtering of gridded data. *Journal of Open Source Software*, 7(70).
- 1005 Loose, N., Bachman, S., Grooms, I., & Jansen, M. (2023). Diagnosing scale-
 1006 dependent energy cycles in a high-resolution isopycnal ocean model. *Journal of*
 1007 *Physical Oceanography*, 53(1), 157–176.
- 1008 Loose, N., Marques, G. M., Adcroft, A., Bachman, S., Griffies, S. M., Grooms, I., ...
 1009 Jansen, M. F. (2023). Comparing two parameterizations for the restratification
 1010 effect of mesoscale eddies in an isopycnal ocean model. *Journal of Advances in*
 1011 *Modeling Earth Systems*, 15(12), e2022MS003518.
- 1012 Mak, J., Maddison, J. R., Marshall, D. P., Ruan, X., Wang, Y., & Yeow, L.
 1013 (2023). Scale-awareness in an eddy energy constrained mesoscale eddy pa-
 1014 rameterization. *Journal of Advances in Modeling Earth Systems*, 15(12),
 1015 e2023MS003886.
- 1016 Marshall, D. P., Maddison, J. R., & Berloff, P. S. (2012). A framework for pa-
 1017 rameterizing eddy potential vorticity fluxes. *Journal of Physical Oceanography*,
 1018 42(4), 539–557.
- 1019 McDougall, T. J., & McIntosh, P. C. (2001). The temporal-residual-mean veloci-
 1020 ty. part ii: Isopycnal interpretation and the tracer and momentum equations.
 1021 *Journal of Physical Oceanography*, 31(5), 1222–1246.
- 1022 Moser, R. D., Haering, S. W., & Yalla, G. R. (2021). Statistical properties of
 1023 subgrid-scale turbulence models. *Annual Review of Fluid Mechanics*, 53,
 1024 255–286.
- 1025 Perezhugin, P., Adcroft, A., & Zanna, L. (submitted). Generalizable neural-network
 1026 parameterization of mesoscale eddies in idealized and global ocean models.
 1027 *Geophysical Research Letters*.
- 1028 Perezhugin, P., Zhang, C., Adcroft, A., Fernandez-Granda, C., & Zanna, L.
 1029 (2024). A stable implementation of a data-driven scale-aware mesoscale pa-
 1030 rameterization. *Journal of Advances in Modeling Earth Systems*, 16(10),
 1031 e2023MS004104.
- 1032 Prakash, A., Jansen, K. E., & Evans, J. A. (2022). Invariant data-driven subgrid
 1033 stress modeling in the strain-rate eigenframe for large eddy simulation. *Com-*
 1034 *puter Methods in Applied Mechanics and Engineering*, 399, 115457.
- 1035 Ramadhan, A., Marshall, J., Souza, A., Wagner, G. L., Ponnampati, M., & Rack-
 1036 auckas, C. (2020). Capturing missing physics in climate model parameteriza-
 1037 tions using neural differential equations. *arXiv preprint arXiv:2010.12559*.
- 1038 Ross, A., Li, Z., Perezhugin, P., Fernandez-Granda, C., & Zanna, L. (2023).
 1039 Benchmarking of machine learning ocean subgrid parameterizations in an
 1040 idealized model. *Journal of Advances in Modeling Earth Systems*, 15(1),

- 1041 e2022MS003258.
- 1042 Sagaut, P. (2005). *Large eddy simulation for incompressible flows: an introduction*.
 1043 Springer Science & Business Media.
- 1044 Sane, A., Reichl, B. G., Adcroft, A., & Zanna, L. (2023). Parameterizing vertical
 1045 mixing coefficients in the ocean surface boundary layer using neural networks.
 1046 *arXiv preprint arXiv:2306.09045*.
- 1047 Smith, K. S. (2007). The geography of linear baroclinic instability in earth’s oceans.
 1048 *Journal of Marine Research*, 65(5), 655–683.
- 1049 Smith, K. S., & Vallis, G. K. (2002). The scales and equilibration of midocean ed-
 1050 dies: Forced–dissipative flow. *Journal of Physical Oceanography*, 32(6), 1699–
 1051 1720.
- 1052 Smith, R. D., & Gent, P. R. (2004). Anisotropic gent–mcwilliams parameterization
 1053 for ocean models. *Journal of Physical Oceanography*, 34(11), 2541–2564.
- 1054 Srinivasan, K., Chekroun, M. D., & McWilliams, J. C. (2023). Turbulence closure
 1055 with small, local neural networks: Forced two-dimensional and β -plane flows.
 1056 *arXiv preprint arXiv:2304.05029*.
- 1057 Tulloch, R., Marshall, J., Hill, C., & Smith, K. S. (2011). Scales, growth rates,
 1058 and spectral fluxes of baroclinic instability in the ocean. *Journal of Physical*
 1059 *Oceanography*, 41(6), 1057–1076.
- 1060 Vallis, G. K. (2017). *Atmospheric and oceanic fluid dynamics*. Cambridge University
 1061 Press.
- 1062 Visbeck, M., Marshall, J., Haine, T., & Spall, M. (1997). Specification of eddy trans-
 1063 fer coefficients in coarse-resolution ocean circulation models. *Journal of physical*
 1064 *oceanography*, 27(3), 381–402.
- 1065 Yuval, J., O’Gorman, P. A., & Hill, C. N. (2021). Use of neural networks for stable,
 1066 accurate and physically consistent parameterization of subgrid atmospheric
 1067 processes with good performance at reduced precision. *Geophysical Research*
 1068 *Letters*, 48(6), e2020GL091363.
- 1069 Zanna, L., & Bolton, T. (2020). Data-driven equation discovery of ocean mesoscale
 1070 closures. *Geophysical Research Letters*, 47(17), e2020GL088376.
- 1071 Zhang, C., Perezhugin, P., Gultekin, C., Adcroft, A., Fernandez-Granda, C., &
 1072 Zanna, L. (2023). Implementation and evaluation of a machine learned
 1073 mesoscale eddy parameterization into a numerical ocean circulation model.
 1074 *Journal of Advances in Modeling Earth Systems*, 15(10), e2023MS003697.