

1 **Uncertainty quantification of ocean parameterizations:**
2 **application to the K-Profile-Parameterization for**
3 **penetrative convection**

4 **A. N. Souza¹, G. L. Wagner¹, A. Ramadhan¹, B. Allen¹, V. Churavy¹, J.**
5 **Schloss¹, J. Campin¹, C. Hill¹, A. Edelman¹, J. Marshall¹, G. Flierl¹, R.**
6 **Ferrari¹**

7 ¹Massachusetts Institute of Technology, 77 Massachusetts Ave, Cambridge, MA 02139, United States

8 **Key Points:**

- 9 • A Bayesian methodology can be applied to turbulence parameterizations to probe
10 parameterizations and better understand their biases and uncertainties.
11 • We can train parameterizations to match large eddy simulations.
12 • We can better understand the physics of parameterizations by applying a Bayesian
13 methodology.

Corresponding author: A. N. Souza, andrenogueirasouza@gmail.com

14 **Abstract**

15 Parameterizations of unresolved turbulent processes in the ocean compromise the fidelity
 16 of large-scale ocean models used in climate change projections. In this work we use a Bayesian
 17 approach for evaluating and developing turbulence parameterizations by comparing pa-
 18 rameterized models with observations or high-fidelity numerical simulations. The method
 19 obtains optimal parameter values, correlations, sensitivities, and more generally, likely
 20 distributions of uncertain parameters. We demonstrate the approach by estimating the
 21 uncertainty of parameters in the popular ‘K-profile parameterization’, using an ensem-
 22 ble of large eddy simulations of turbulent penetrative convection in the ocean surface bound-
 23 ary layer. We uncover structural deficiencies and discuss their cause. We conclude by
 24 discussing the applicability of the approach to Earth system models.

25 **Plain Language Summary**

26 Climate projections continue to be marred by large uncertainties, which originate
 27 in the poor representation of physical processes that occur at scales too small for climate
 28 models to properly simulate them, like clouds in the atmosphere and turbulent swirls
 29 in the ocean. We propose to develop more accurate representations of small physical ocean
 30 processes (parameterizations) trained with high resolution numerical simulations of small
 31 ocean patches. A Bayesian methodology is used to calibrate the parameterizations with
 32 the high resolution numerical simulations, to assess their fidelity and to identify improve-
 33 ments. Most importantly this approach provides estimates of the uncertainties in the pa-
 34 rameterizations which can then be used to quantify uncertainties of climate models. While
 35 the approach is illustrated for a parameterization of ocean turbulence, it can be applied
 36 to any parameterization in climate models.

37 **1 Introduction**

38 The ocean components of Earth system models are complex systems that couple
 39 the resolved ocean circulation with a myriad of unresolved, parameterized and impor-
 40 tant physical processes. Parameterizations of unresolved physical processes often involve
 41 many uncertain parameters which are used to tune the model in an attempt to obtain
 42 a desired outcome (Hourdin et al., 2017). Moreover, each component, whether resolved
 43 or parameterized, interacts with all the others in nonlinear ways that lead to complex
 44 behavior which is sometimes difficult to understand and characterize.

45 Upper ocean turbulent mixing is a key parameterized process in ocean circulation
 46 models. The detailed fluid dynamics of upper ocean turbulent mixing are highly com-
 47 plex, involving surface boundary layer turbulence driven by buoyancy loss or winds, bot-
 48 tom boundary layer turbulence, lateral mixing due to baroclinic effects, and so forth. How-
 49 ever, at least in principle, the governing fluid dynamical equations are known. The prob-
 50 lem is that the computational resources required to resolve them and, at the same time,
 51 the global scale circulation, are not available and will not be for the foreseeable future
 52 (Schneider, Teixeira, et al., 2017).

53 A goal of this paper is to outline and illustrate a Bayesian framework to assess and
 54 improve parameterizations. We present a way forward which employs an ensemble of mix-
 55 ing process resolving simulations to train a chosen parameterization. The core idea is
 56 that the parameterization must represent the collective effect of sub-grid scale physics
 57 faithfully for all relevant relevant external forcings and mean climate states. This con-
 58 trasts approaches that attempt to diagnose parameters directly from high-resolution sim-
 59 ulations or to estimate values that perform well only in a particular experiment. It should
 60 be noted, however, that by restricting ourselves to understanding parameterizations in
 61 the context of sub-grid scale physics, we may miss out on important interactions with
 62 the rest of the climate system, e.g., the interaction of resolved lateral fluxes from the global

63 ocean circulation with parameterized turbulent vertical mixing in the ocean. Neverthe-
 64 less, studying one subgrid-scale process at a time is not an exercise in futility since it is
 65 a necessary first step to optimize a parameterization before considering the interactions
 66 with all other components of the full system.

67 We take a Bayesian perspective in our optimization of parameterizations. There
 68 are many ways in which a Bayesian framework can be used. Here we will explore one
 69 particular approach: characterizing the parameters of a parameterization via probabil-
 70 ity distributions. Thus, we will go beyond finding a point estimate for parameters. These
 71 probability distributions capture the notion of uncertainty and nonlinear correlations be-
 72 tween parameters. Furthermore, they can then be used as prior distributions for param-
 73 eter sensitivity studies in full climate models. This partially addresses a present deficiency
 74 in the current approach used to tune parameters in climate models. “Manual” tuning
 75 is done to obtain agreement between models and observations (Hourdin et al., 2017). Since
 76 parameters are often correlated, a parameter may be tuned to offset biases introduced
 77 by another parameter, resulting in parameterizations that no longer respect the subgrid-
 78 scale physics. The Bayesian framework automates parameter search in a way that en-
 79 sures it respects the underlying physics of a parameter. The calibration of parameter-
 80 ization schemes in climate models has the potential to reduce biases as well as quantify
 81 the uncertainty of key climate variables, such as ocean heat content or climate sensitiv-
 82 ity; however, innovation is required to make the Bayesian method practical and compu-
 83 tationally feasible in the global model. One step towards this is to calculate prior dis-
 84 tributions for parameters in a simplified setting, such as the local studies performed here,
 85 and then use computationally efficient methods for obtaining posterior distributions in
 86 the global climate model such as those proposed in (Schneider, Lan, et al., 2017; Albers
 87 et al., 2019; Cleary et al., 2020).

88 The focus here is to calculate prior distributions for parameters in ocean climate
 89 models. We do so by matching parameterizations to large eddy simulations, a philoso-
 90 phy similar in spirit to that which has been done in the atmospheric context for cloud
 91 parameterizations (Golaz et al., 2007). To make our discussion concrete we focus on the
 92 representation of convectively-driven turbulence in the upper ocean.

93 Our paper is organized as follows: In section 2 we describe the physical scenario
 94 in which we run our Large Eddy Simulations (LES) and parameterization. In section 3
 95 we introduce Bayesian parameter estimation for the parameters in the K-Profile Param-
 96 eterization (KPP) and perform the parameter estimation in the regime described by sec-
 97 tion 2. Finally, we end with a discussion in section 4.

98 **2 Large eddy simulations and K-profile parameterization of penetra-** 99 **tive convection**

100 During the onset of winter at high latitudes, cooling at the ocean surface gener-
 101 ates convective plumes that descend and mix the ocean surface boundary layer, see Marshall
 102 and Schott (1999) for a review. Near-surface mixing by convection generates a surface
 103 layer of uniform temperature and salinity called the ‘mixed layer’ which can reach depths
 104 of hundreds of meters.

105 At the base of the mixed layer, convective plumes penetrate further into a strongly-
 106 stratified region called the ‘entrainment layer’, where plume-driven turbulent mixing be-
 107 tween the mixed layer and the ocean interior further cools the boundary layer. This pro-
 108 cess, in which the surface layer is cooled both at the surface and by turbulent mixing in
 109 the entrainment layer, is called penetrative convection. Penetrative convection is a cru-
 110 cial oceanic process for storing heat and carbon as well as setting the density structure
 111 of the deep ocean. Parameterizations of ocean surface boundary layer mixing must de-
 112 scribe penetrative convection accurately. In this paper we evaluate the accuracy of the

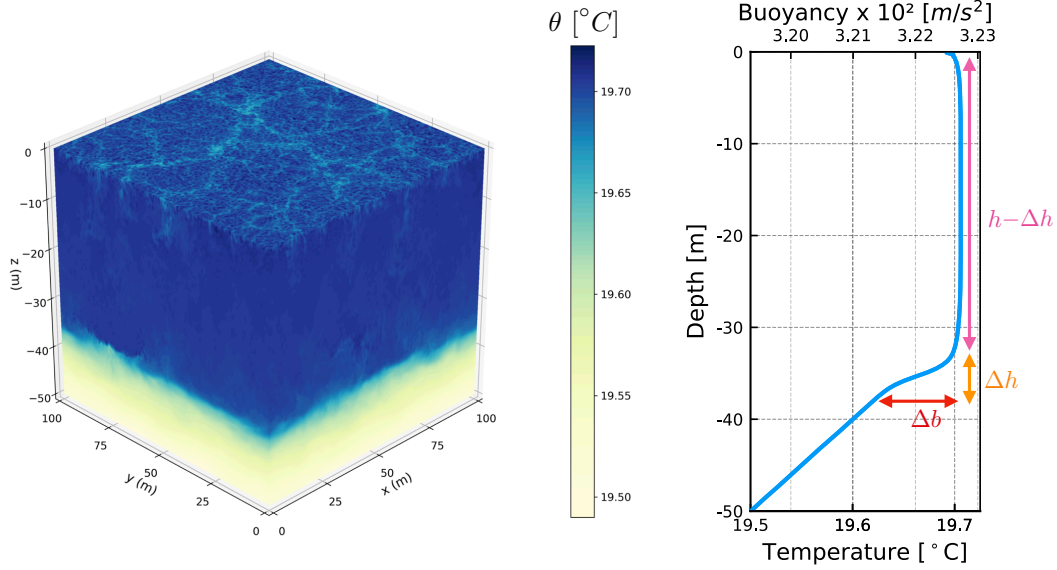


Figure 1. A 3D simulation of the LES model of the Boussinesq equations and its horizontal average at $t = 2$ days.

113 K -profile parameterization (Large et al., 1994) against large eddy simulations (LES) of
 114 idealized penetrative convection scenarios.

115 2.1 Idealized penetrative convection scenario

Our idealized scenarios impose a constant surface cooling $Q_h > 0$ to a resting, linearly stratified boundary layer with the initial state

$$\mathbf{u}|_{t=0} = 0 \text{ and } b|_{t=0} = N^2 z, \quad (1)$$

where $\mathbf{u} = (u, v, w)$ is the resolved velocity field simulated by LES, b is buoyancy, and N^2 is the initial vertical buoyancy gradient. The surface buoyancy flux Q_b is related to the imposed surface cooling Q_h , which has units W m^{-2} , via

$$Q_b = \frac{\alpha g}{\rho_{\text{ref}} c_p} Q_h, \quad (2)$$

116 where $\alpha = 2 \times 10^{-4} (\text{°C})^{-1}$ is the thermal expansion coefficient, $g = 9.81 \text{ m s}^{-2}$ is grav-
 117 itational acceleration, $\rho_{\text{ref}} = 1035 \text{ kg m}^{-3}$ is a reference density, and $c_p = 3993 \text{ J/(kg °C)}$
 118 is the specific heat capacity. Our software and formulation of the large eddy simulations
 119 is discussed in section Appendix A.

Output of a large eddy simulation of turbulent penetrative convection in a domain $L_x = L_y = L_z = 100$ meters is in Figure 1. The left panel in Figure 1 visualizes the three-dimensional temperature field $\theta = \theta_0 + b/\alpha g$ associated with the buoyancy b , where $\theta_0 = 20^\circ\text{C}$ is the surface temperature at $z = 0$. The right panel of Figure 1 shows the horizontally averaged buoyancy profile

$$\bar{b}(z, t) \equiv \frac{1}{L_x L_y} \int_0^{L_x} \int_0^{L_y} b(x, y, z, t) dx dy. \quad (3)$$

120 The visualization reveals the two-part boundary layer produced by penetrative con-
 121 vection: close to the surface, cold and dense convective plumes organized by surface cool-

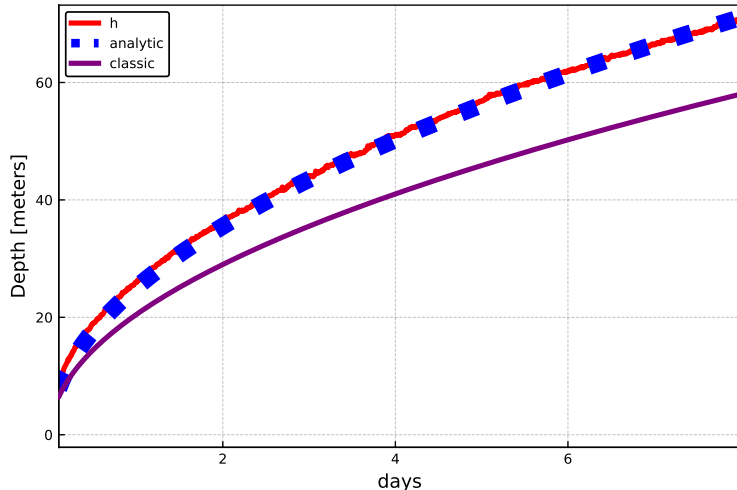


Figure 2. Mixed layer depth and its evolution in time after initial transients. The blue squares are the analytic scaling 4, the red line is an estimate of the boundary layer depth directly from the LES (described in the text), and the purple line is the classic scaling which ignores the entrainment layer 8.

122 ing sink and mix ambient fluid, producing a well-mixed layer that deepens in time. Be-
 123 low the mixed layer, the momentum carried by sinking convective plumes leads them to
 124 overshoot their level of neutral buoyancy (nominally, the depth of the mixed layer), ‘pen-
 125 etrating’ the stably stratified region below the surface mixed layer and generating the
 126 strongly stratified entrainment layer. The total depth of the boundary layer is h and in-
 127 cludes the mixed layer and the entrainment layer of thickness Δh . Turbulent fluxes are
 128 negligible below $z = -h$ for our purposes.

In figure 2 we show the evolution of $h(t)$ defined as the first depth from the bot-
 tom where the stratification is equal to a weighted average of the maximum stratifica-
 tion and the initial stratification¹. The dotted line confirms that the evolution after an
 initial transient is best fit by the formula,

$$h \simeq \sqrt{3.0 \frac{Q_b}{N^2} t}, \quad (4)$$

129 where N^2 is the initial stratification.

This result is easily explained by considering the horizontally averaged buoyancy
 equation,

$$\partial_t \bar{b} = -\partial_z (\bar{wb} + \overline{q^{(z)}}), \quad (5)$$

where \bar{b} is the horizontally averaged buoyancy, \bar{wb} is the horizontally averaged vertical
 advective flux and $\overline{q^{(z)}}$ is the horizontally averaged vertical diffusive flux. Integrating the
 equation in time between $t' = 0$ and some later time $t' = t$, and in the vertical be-
 tween the surface, where $q^{(z)} = -Q_b$, and the base of the entrainment layer where all

¹ The weights are $2/3$ for the initial stratification N^2 and $1/3$ for the maximum stratification N_m^2 so
 that h satisfies $\partial_z \bar{b}(-h) = 2N_b^2/3 + N_m^2/3$. This guarantees that h is a depth where the local stratification
 lies between the background stratification and the maximum stratification since it is defined as the *first*
 depth starting from the bottom that satisfies such a criteria.

turbulent fluxes vanish, one finds,

$$\int_{-h}^0 [\bar{b}(z, t) dz - \bar{b}(z, 0)] dz = -Q_b t. \quad (6)$$

Substituting $\bar{b}(z, 0) = b_0 + N^2(z+h)$ and $\bar{b}(z, t) = b_0 + \Delta b$, which is an appropriate approximation of the profile shown in Fig. 1b except at very early times in the simulation, yields

$$\frac{1}{2}N^2h^2 - h\Delta b = Q_b t. \quad (7)$$

The first term on the left of equation 7 describes boundary layer deepening due to buoyancy extraction at the surface, while the second term corresponds to the further cooling caused by turbulent mixing in the entrainment layer. Ignoring turbulent mixing in the entrainment layer yields the deepening rate

$$h = \sqrt{2.0 \frac{Q_b}{N^2} t}, \quad (8)$$

130 which differs by roughly 20% from the best fit expression 4 due to the effects of turbu-
 131 lent mixing in the entrainment layer. The scaling in equation 8 is the deepening rate as-
 132 sociated with a convective adjustment parameterization. The K -profile parameteriza-
 133 tion of penetrative convection, on the other hand, introduces a model for entrainment
 134 layer mixing in an attempt to describe equation 4.

135 2.2 The K -profile parameterization of penetrative convection

In penetrative convection in a horizontally-periodic domain, the K -profile parameterization models the evolution of the horizontally averaged temperature profile, $\bar{\theta}(z, t)$, and the boundary layer depth with

$$\partial_t T = -\partial_z F(T, h; \mathbf{C}) \quad (9)$$

$$0 = \mathcal{D}(T, h; \mathbf{C}), \quad (10)$$

136 where T is the temperature profile produced by the K -profile parameterization, h is the
 137 boundary layer depth, $\mathbf{C} = \{C^S, C^N, C^D, C^H\}$ is a set of free parameters for represent-
 138 ing dimensionless proportionality constants following various scaling laws, $F(T, h; \mathbf{C})$ is
 139 a ‘temperature flux function’, and $\mathcal{D}(T, h; \mathbf{C})$ is a nonlinear-integral constraint to deter-
 140 mine the boundary layer depth. We emphasize that the goal of the parameterization is
 141 not limited to just getting the mixed layer depth correct or correctly predicting the jump
 142 in buoyancy. Rather, the goal is to obtain correct heat exchanges with the atmosphere,
 143 entrainment of nutrients for the biology, and flux rates of passive scalars. Thus, it is im-
 144 portant to faithfully capture the dynamics of the entire temperature profile.

The K -profile parameterization (KPP) represents F through the sum of a down-
 gradient flux and a non-local flux term (Large et al., 1994),

$$F = - \underbrace{C^D w_* h \frac{z}{h} \left(1 + \frac{z}{h}\right)^2}_{\equiv K} \partial_z T + \underbrace{C^N Q^\theta \frac{z}{h} \left(1 + \frac{z}{h}\right)^2}_{\equiv \Phi}, \quad (11)$$

145 for $-h \leq z \leq 0$ and 0 otherwise. Here $w_* = (Q_b h)^{1/3}$ is the convective turbulent ve-
 146 locity scale, h is the boundary layer depth, $\frac{z}{h} \left(1 + \frac{z}{h}\right)^2$ is the ‘ K -profile’ shape function—
 147 K is the namesake downgradient diffusivity of the K -profile parameterization—and Φ
 148 is a ‘non-local’ flux term that models convective boundary layer fluxes not described by
 149 downgradient diffusion.

In penetrative convection, the KPP model estimates the boundary layer depth h
 with an implicit nonlinear equation. To motivate the functional form of this criteria, first

see figure 1 for reference. The jump in buoyancy, Δb , is the difference between the buoyancy in the mixed layer and the base of the entrainment region. Equivalently we can write $\Delta b = N_e^2 \Delta h$, where N_e the stratification in the entrainment region. From plume theory, see section Appendix B, we obtain $\Delta h \propto w_\star / N_e$ so that

$$C^H = \frac{\Delta b}{w_\star N_e} \quad (12)$$

for some universal proportionality constant C^H , which we call the ‘mixing depth’ parameter. KPP posits that the mixed layer depth h is the first such depth from the surface at which equation 12 holds. For numerical stability reasons equation 12 is generally formulated as

$$C^H = \frac{h \left[\frac{1}{C^S h} \int_{-C^S h}^0 B(z) dz - B(-h) \right]}{(hQ_b)^{1/3} h \sqrt{\max[0, \partial_z B(-h)] + 10^{-11} \text{m}^2 \text{s}^{-2}}}, \quad (13)$$

150 where $B = \alpha g T$. The numerator approximates the jump in buoyancy times the mixed
 151 layer depth, $h \Delta b$. The term, $\frac{1}{C^S h} \int_{-C^S h}^0 B(z) dz$, serves as an approximation to the buoy-
 152 ancy in the mixed layer. The denominator evaluates the product $h w_\star N_e$ with $w_\star = (h Q_b)^{1/3}$
 153 and $N_e = \sqrt{\max[0, \partial_z B(-h)]}$ at a given depth, while adding a dimensional term 10^{-11}
 154 to prevent division by zero. In section Appendix B we go in further detail about the rati-
 155 onale behind the implicit equation for the boundary layer depth, equation 13, for the
 156 case of penetrative convection.

157 The mixing depth parameter, C^H , is often referred to as the critical bulk Richard-
 158 son number in the KPP literature (Large et al., 1994), because in mechanically forced
 159 turbulence, the denominator is replaced by the mean shear squared times h . In pene-
 160 trative convection there is no mean shear and C^H is no longer related to a bulk Richard-
 161 son number.

The K -profile parameterization for penetrative convection has four free paramete-
 162 rs: the surface layer fraction C^S , the flux scalings C^N and C^D in equation 11, and the
 163 mixing depth parameter C^H in equation 13. Their default values, reported in (Large et
 al., 1994), are

$$(C^S, C^N, C^D, C^H) = (0.1, 6.33, 0.77, 0.95). \quad (14)$$

162 Our objective is to calibrate the free parameters $\mathbf{C} = (C^S, C^N, C^D, C^H)$ by compar-
 163 ing KPP temperature profiles $T(z, t; \mathbf{C})$ with the LES output $\bar{\theta}(z, t)$.

164 3 Model Calibration

We outline a Bayesian method for optimizing and estimating the uncertainty of the
 165 four free parameters through a comparison of solutions $T(z, t; \mathbf{C})$ to equation 9 to the
 166 output $\bar{\theta}(z, t)$ of our large eddy simulations. For this we define a loss function by

$$\mathcal{L}(\mathbf{C}) = \max_{t \in [t_1, t_2]} \left\{ \frac{1}{L_z} \int_{-L_z}^0 [T(z, t; \mathbf{C}) - \bar{\theta}(z, t)]^2 dz \right\}. \quad (15)$$

165 We choose the square error in space to reduce the sensitivity to vertical fluctuations in
 166 the temperature profile. In time we take the maximum value of the squared error to guar-
 167 antee that the temperature profile never deviates too far from the LES simulation at each
 168 instant.

169 Notably we do not use the boundary layer depth in the definition of the loss func-
 170 tion. Firstly, it should be stressed that getting the entire temperature profile correct is
 171 a more stringent requirement and would also imply a correct boundary (and mixed) layer
 172 depth. We prefer not to use a boundary layer depth directly because it leads to noisy

173 loss functions and depends too much on the precise definition used. In the literature, there
 174 are several different kinds of “depth” parameters based on, for example, the KPP def-
 175 inition as per equation 13, the location of the minimum buoyancy flux, the location of
 176 the maximum temperature gradient, or the first depth at which the temperature decreases
 177 by some ΔT of the surface value (Kara et al., 2000; Van Roekel et al., 2018). It is sim-
 178 pler (albeit more ambitious) to target the entire temperature profile. We prefer not to
 179 use the horizontally averaged temperature fluxes or gradients for practical reasons. Fluxes
 180 tend to be noisier than the horizontally averaged temperature profile and one would have
 181 to apply a smoothing filter. In summary, these other metrics introduce additional sources
 182 of systematic bias for little gain in the present circumstance.

183 A natural way to extend the definition of loss functions in order to take into ac-
 184 count parameter sensitivities is to define probability distributions for parameters. Sim-
 185 ilar to how the functional form of the loss function is critical to the estimation of opti-
 186 mal parameters, the functional form of a probability density is critical for estimating the
 187 uncertainties of a parameter. A probability distribution quantifies what we mean by “good”
 188 or “bad” parameter choices, (similar to a loss function), but in terms of uncertainties and
 189 likelihoods. It is often the case that one has a good feel for how to define meaningful loss
 190 functions, but less so for probability distributions. Here we report our choices, but in sec-
 191 tion Appendix C we provide guidance on criteria to be used when constructing as well
 192 as sampling from the probability distribution. It is worth keeping in mind that, just like
 193 loss functions, the true test is “after-the-fact”; we inspect results and confirm that they
 194 indeed correspond to our intuition. Just like the definition of a loss function implicitly
 195 determines a choice of optimal parameters, a choice of probability distribution implic-
 196 itly determines parameter sensitivities². Both are arbitrary, but that does not mean that
 197 loss functions or parameter sensitivities are meaningless.

We adopt the same definition as in (Schneider, Lan, et al., 2017) for the probabil-
 ity distribution:

$$\rho(\mathbf{C}) \propto \rho^0(\mathbf{C}) \exp\left(-\frac{\mathcal{L}(\mathbf{C})}{\mathcal{L}_0}\right) \quad (16)$$

198 where ρ^0 is the prior distribution of the parameter values, \mathcal{L} is a loss function, and $\mathcal{L}_0 >$
 199 0 is a *hyperparameter*³.

The loss function \mathcal{L} has dimensions and the parameter \mathcal{L}_0 makes the quantity in
 the exponent dimensionless. \mathcal{L}_0 could have been absorbed into the loss function, but it
 has a probabilistic interpretation that is worth emphasizing. We chose the parameter
 \mathcal{L}_0 as the minimum of the loss function $\mathcal{L}(\mathbf{C})$ —the minimum is found using a modified⁴
 simulated annealing procedure to compute the minimum of \mathcal{L} (Kirkpatrick et al., 1983).
 With this choice the likelihood of any other parameter choice, say \mathbf{C}^1 , is determined by
 the amount by which it increases the minimum of the loss function, i.e.,

$$\rho(\mathbf{C}^1)/\rho(\mathbf{C}^*) = \exp\left(\frac{\mathcal{L}^0 - \mathcal{L}(\mathbf{C}^1)}{\mathcal{L}^0}\right), \quad (17)$$

200 where \mathbf{C}^* denotes the optimal parameter choice with $\mathcal{L}^0 = \mathcal{L}(\mathbf{C}^*)$. For example, if the
 201 choice \mathbf{C}^1 increases the minimum of the loss function by a factor of two, i.e. $\mathcal{L}(\mathbf{C}^1) =$
 202 $2\mathcal{L}^0$, then it is $1/e$ less likely.

² Parameter sensitivities are inversely related to parameter uncertainties. A more sensitive parameter
 is one that produces larger changes to the loss function. In the context of this paper a more uncertain
 parameter is one that produces small changes to the loss function. See section ?? for a simple example.

³ A hyperparameter is a parameter associated with the probability distribution as opposed to a parame-
 ter in the parameterization.

⁴ The main difference is that we take the minimum “artificial temperature associated with the simu-
 lated annealing procedure” to be the best known minimum of the loss function \mathcal{L} rather than 0.

203 Once \mathcal{L}_0 is determined, we use the Random Walk Markov Chain Monte Carlo (RW-
 204 MCMC) algorithm (Metropolis et al., 1953), described further in section C2, to sample
 205 the probability distribution.

It is worth mentioning that equation 16 is the continuous analogue of Bayes formula

$$\mathbb{P}(\mathbf{C}|\text{data}) \propto \mathbb{P}(\mathbf{C})\mathbb{P}(\text{data}|\mathbf{C}) \quad (18)$$

206 where \mathbb{P} is a probability distribution. In our context we interpret the formula as follows:
 207 We update our prior belief of the distribution of parameters \mathbf{C} based on the data (in this
 208 case the LES experiment). $\mathbb{P}(\mathbf{C})$ is our prior probability for the parameters \mathbf{C} , while $\mathbb{P}(\text{data}|\mathbf{C})$
 209 is the probability that the parameter choices \mathbf{C} explain the data. Choosing to model pa-
 210 rameters as probability distributions has the consequence that the output of the param-
 211 eterization is also inherently probabilistic. In particular, the output of KPP will no longer
 212 be just a point estimate for temperature at each depth and each moment in time, but
 213 rather a probability distribution.

For all the uncertainty quantification that follows, we use resolution and timesteps
 typical of state of the art ocean models used for climate studies: a resolution of 100 m/16 =
 6.25 m and a timestep of ten minutes. The temporal window used to compute the loss
 function is from $t_1 = 0.25$ days to the final simulation day. We apply the Bayesian pa-
 rameter estimation procedure to KPP using data from one LES simulation in section 3.1
 and from multiple LES simulations using different initial stratifications in section 3.2.
 We use a uniform prior for the parameters in KPP over the following ranges:

$$0 \leq C^S \leq 1, \quad 0 \leq C^N \leq 8, \quad 0 \leq C^D \leq 6, \quad \text{and} \quad 0 \leq C^H \leq 5. \quad (19)$$

214 The surface layer fraction C^S , being a fraction, must stay between zero and one. The
 215 other parameter limits were chosen to correspond to “reasonable” ranges around the de-
 216 fault values, equation 14.

217 3.1 Calibration of KPP parameters against one LES simulation

218 In this section we apply the Bayesian calibration method to the LES simulation
 219 of penetrative convection described in section 2.1 and quantify uncertainties in param-
 220 eters of KPP, section 2.2. The horizontal averages from the LES simulations are com-
 221 pared with predictions from solutions of the KPP diffusion scheme. The boundary and
 222 initial conditions for KPP are taken to be the same as those for the LES simulation, i.e.,
 223 100 W/m² cooling at the top, $\partial_z T = 0.01^\circ\text{C m}^{-1}$ at the bottom, and an initial pro-
 224 file $T_p(z, 0) = 20^\circ\text{C} + 0.01^\circ\text{C m}^{-1}z$.

We use the RW-MCMC algorithm with 10^6 iterations to sample the probability dis-
 tributions of the four KPP parameters (C^S, C^N, C^D, C^H). This lead to roughly 10^4 sta-
 tistically independent samples as estimated using an autocorrelation length, see Sokal
 (1997). The RW-MCMC algorithm generates the entire four dimensional PDF, equation
 16, but visualizing this object is challenging. Instead we look at the marginal distribu-
 tions, e.g.,

$$\rho_M(C^H) \equiv \iiint \rho(\mathbf{C}) \, dC^S dC^D dC^N, \quad (20)$$

225 and similarly for the other parameters. (Constructing the marginal distributions only
 226 requires constructing histograms of the trajectories generated by the RW-MCMC algo-
 227 rithm.) Parameter correlations are washed away by focusing on marginal distributions.
 228 Nevertheless, marginal distributions give the range of parameter values that yield little
 229 change to the loss function and are shown in figure 3. The marginal distribution of the
 230 mixing depth parameter C^H is much more compact than that of the other three param-
 231 eters suggesting that it is the most sensitive parameter. The mixing depth parameter’s

232 importance stems from its control over both the buoyancy jump across the entrainment
 233 layer and the rate-of-deepening of the boundary layer. (Once again it may be useful to
 234 remember that C^H is often referred to the bulk Richardson number in the KPP liter-
 235 ature, even though it take a different meaning in convective simulations.) The param-
 236 eters C^K and C^N set the magnitude of the local and nonlocal fluxes and their specific
 237 value is not too important as long as they are large enough to maintain a well-mixed layer.
 238 The value of the regularization C^S is quite irrelevant.

239 The parameter distribution can be used to choose an optimal set of KPP param-
 240 eters. Of the many choices, we choose the most probably value of the four dimensional
 241 probability distribution, the mode, because they minimize the loss function as explained
 242 in section Appendix C. (These values do not necessarily correspond to the individual modes
 243 of the marginal distributions. For example C^H is set to ≈ 2.0 rather than 1.5.) In fig-
 244 ure 4a we show the area averaged temperature profile after 8 days from the LES sim-
 245 ulation (continuous line) and the temperature profiles obtained running the KPP param-
 246 eterization with default and optimal parameters (squares and dtots). The optimized tem-
 247 perature profiles are more similar to the LES simulation than the default value especially
 248 in the entrainment region. figure 4b confirms that the square root of the loss function,
 249 the error, grows much faster with the default parameters. The oscillations in the error
 250 are a consequence of the coarseness of the KPP model: only one grid point is being en-
 251 trained at any given moment.

252 The improvement in boundary layer depth through optimization of the paramet-
 253 ers is about 10%, or 10 m over 8 days. As we discussed in section 2.1, the rate of deep-
 254 ening can be predicted analytically within 20% by simply integrating over time and depth
 255 the buoyancy budget and assuming that the boundary layer is well mixed everywhere,
 256 i.e. ignoring the development of enhanced stratification within an entrainment layer at
 257 the base of the mixed layer. KPP improves on this prediction by including a paramet-
 258 erization for the entrainment layer. The default KPP parameters contribute a 10% im-
 259 provement on the no entrainment layer prediction, and the optimized parameters con-
 260 tribute another 10%. While these may seem like modest improvements, they can result
 261 into large biases in boundary layer depth when integrated over a few months of cooling
 262 in winter rather than just 8 days. We will return to this point in the next section when
 263 we discuss structural deficiencies in the KPP formulation.

264 The probability distributions of the parameters can be used to predict the prob-
 265 ability distributions of all variables, for example temperature at each depth and time,
 266 predicted by KPP. To do this, we subsample the 10^6 parameter values down to 10^4 and
 267 evolve KPP forward in time for each set of parameter choices. We construct histograms
 268 for the temperature field at the final time for each location in space individually. We then
 269 stack these histograms to create a visual representation of the model uncertainty. This
 270 uncertainty quantifies the sensitivity of the parameterization with respect to paramet-
 271 er perturbations as defined by the parameter distributions.

272 The histogram of temperature profiles at time $t = 8$ days as calculated by both
 273 our prior distribution (uniform distribution) and the posterior distribution (as obtained
 274 from the RW-MCMC algorithm) is visualized in figure 5. We see that there is a reduc-
 275 tion of the uncertainty in the temperature profile upon taking into account information
 276 gained from the LES simulation. The salient features of the posterior distribution tem-
 277 perature uncertainty are

- 278 1. 0-10 meter depth: There is some uncertainty associated with the vertical profile
 279 of temperature close to the surface.
- 280 2. 20-60 meter depth: The mean profile of temperature in the mixed layer is very well
 281 predicted by KPP.
- 282 3. 60-70 meter depth: The entrainment region contains the largest uncertainties.

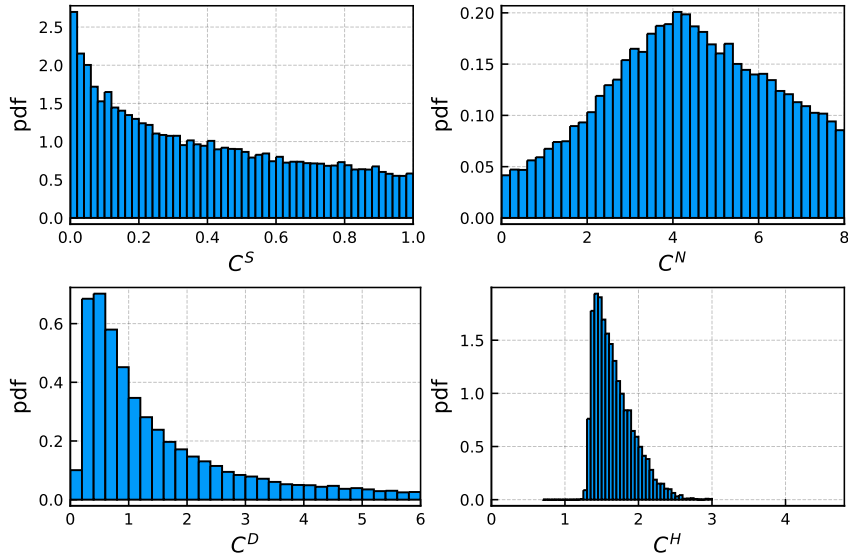


Figure 3. Parameter marginal posterior probability distributions. Marginal probability correspond to parameters parameters: C^S Surface Layer Fraction, C^N nonlocal diffusivity amplitude, C^D diffusivity amplitude, C^H mixing depth parameter. The probability distributions capture the notion of what parameter values are “good” and which ones are “bad”. For example, in the pdf for C^H we see that a value of 2.5 is probable but a value of 5 would be not be. This intuitively corresponds to saying that a value of 2.5 would be a “reasonable” choice whereas 5 would be “unreasonable”. The width of the C^S and C^N parameters suggest that KPP is quite insensitive to their values. A similar consequence holds for C^D , but there also seems to be a preference for values around one.

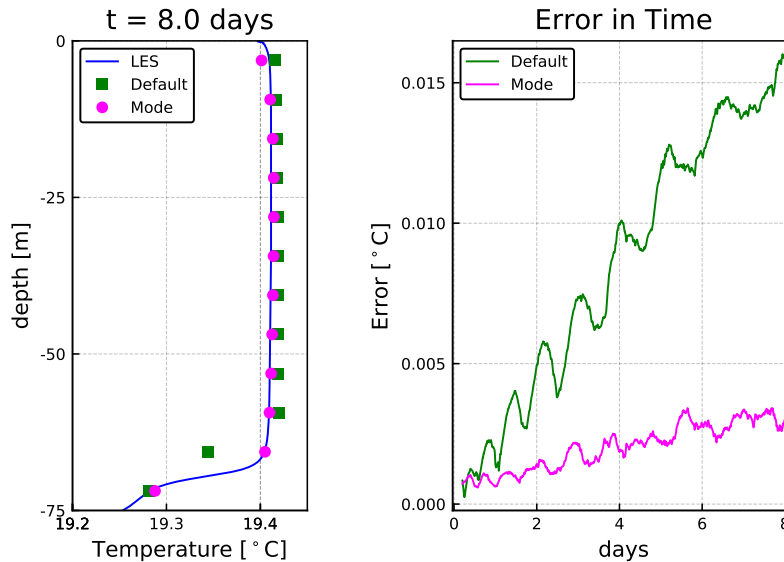


Figure 4. KPP and horizontally averaged LES temperature profiles for different point estimates of parameters at $t = 8$ days as well as the error in time. In the left plot, the squares correspond to default parameter choices, the circles correspond to the optimized parameterization (using the mode of the probability distribution), and the blue line to the horizontally averaged LES solution, all at time $t = 8$ days. On the right plot we show the instantaneous error at each moment in time. We see that the “optimal” parameter does indeed reduce the bias over the time period. The loss function is the largest square of the error over the time interval.

283 4. 70-100 meter depth: There is virtually no uncertainty. The unstratified region be-
 284 low the boundary layer does not change from its initial value.

285 Now that we have applied the Bayesian methodology to one LES simulation and
 286 explored its implications, we are ready to apply the method to multiple LES simulations
 287 covering different regimes in the following section. We focus on the optimization and un-
 288 certainty quantification of C^H for the remainder of the paper, since it is the most sen-
 289 sitive parameter. In the background, we are estimating *all* parameters.

290 3.2 Calibration of KPP parameters from multiple LES simulation

291 There are many possible directions that one could take at this point. We present
 292 an example of how we can use the methodology to explore bias in the KPP model. To
 293 this end we investigate what happens when we change the initial stratification in pen-
 294 etrative convection simulations. This is an informed decision motivated by recent work
 295 on mixed layer depth biases in the Southern Ocean (DuVivier et al., 2018; Large et al.,
 296 2019). In those studies, KPP failed to simulate deep mixed layer in winters when the sub-
 297 surface summer stratification was strong.

298 We perform 32 large eddy simulations and calculate parameter distributions for each
 299 case. We kept the surface cooling constant at 100 W/m^2 for all regimes, and only var-
 300 ied the initial stratification. The integration time was stopped when the boundary layer
 301 depth filled about 70% of the domain in each simulation. We used 128^3 grid points, \approx

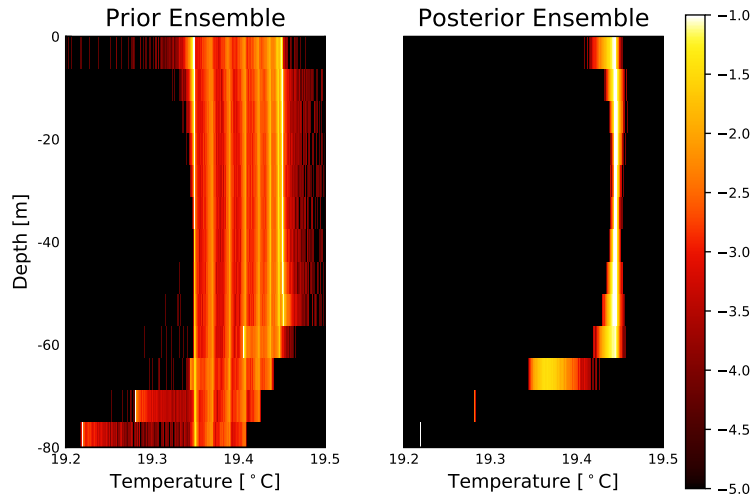


Figure 5. Uncertainty propagation of the temperature profile with respect to the prior and posterior probability distributions. The use of probability distributions for parameters has the consequence that the temperature field is no longer a point estimate, but rather a probability distribution at each moment in space and time. By sampling from the parameter probability distributions and evolving the parameterization forward in time, we obtain a succinct representation of what it means to “fiddle” with parameters. The legend on the right shows what the colors correspond to in terms of the base 10 logarithm of the probability distributions.

302 0.8 meter resolution in each direction⁵. Each one of the probability distributions used
303 10^5 iterations of RW-MCMC, leading to effective sample size on the order of 10^3 .

304 The result, which is visualized in figure 6, shows that the parameter C^H depends
305 on the background stratification, N^2 . The blue dots are the median values of the prob-
306 ability distributions and the stars are the modes (minimum of the loss function). The
307 error bars correspond to 90% probability intervals, meaning that 90% of parameter val-
308 ues fall between the error bars. The default KPP value is plotted as a dashed line for
309 reference.

310 The median values and optimal values increase monotonically with the initial strat-
311 ification value. Given that the parameter is supposed to be dimensionless, this reveals
312 a systematic bias. Furthermore, it exposes *where* the systematic bias comes from: the
313 boundary layer depth criterion in equation 13. No single value of C^H can correctly re-
314 produce the deepening of the boundary layer for all initial stratifications.

The failure of the depth criterion can be understood by going back to the buoy-
ancy budget in equation 7. Using the KPP estimate for the buoyancy jump across the
entrainment layer,

$$\Delta b \equiv \frac{1}{C^S h} \int_{-C^S h}^0 B(z) dz - B(-h), \quad (21)$$

and introducing $N_h^2 \equiv \partial_z B(-h)$ for the stratification at the base of the entrainment
layer to distinguish it from the interior stratification N^2 , we find that the boundary layer

⁵ Although the parameter estimates will vary upon using less resolution, the qualitative trends are expected to be robust.

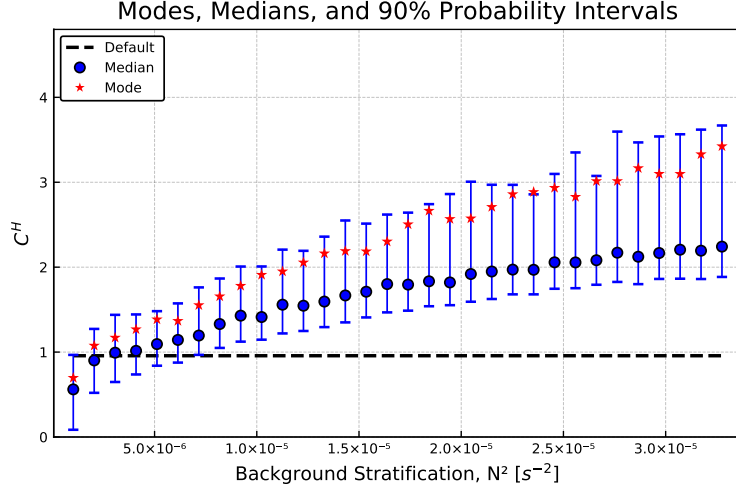


Figure 6. Mixing depth parameter optimized across various background stratification. The dots are the median values, the stars are the mode, and the error bars correspond to 90% probability intervals. The horizontal dashed line is the default value of the mixing depth parameter for reference. Here one can see that there mixing depth parameter when estimated across various regimes produces different results. This is a signature of a systematic bias in the parameterization.

depth criterion, equation 13, implies,

$$h\Delta b \simeq C^H h^{4/3} (Q_b)^{1/3} N_h. \quad (22)$$

Substituting this expression in the buoyancy budget, equation 7, one obtains an implicit equation for the evolution of the boundary layer depth h ,

$$\left(\frac{1}{2} N^2 - C^H (Q_b)^{1/3} h^{-2/3} N_h \right) h^2 \simeq Q_b t. \quad (23)$$

315 The LES simulation described in section 2.1, and many previous studies of penetrative
 316 convection, e.g. (Van Roekel et al., 2018; Deardorff et al., 1980), show that the bound-
 317 ary layer depth grows as \sqrt{t} . N_h would have to scale as $h^{2/3}$ for KPP to correctly pre-
 318 dict that deppening rate of the boundary layer, but this scaling is not observed in the
 319 LES simulations nor supported by theory.

3.3 Modification

From the multi-regime study of the previous section we found that there is no optimal KPP mixing depth parameter C^H that works for arbitrary initial stratification. This prompted us to look for an alternative formulation of the depth criterion which satisfies the well known empirical result that the boundary layer depth deepens at a rate,

$$h \simeq \sqrt{c \frac{Q_b}{N^2} t}, \quad (24)$$

where c is a dimensionless constant found to be close to 3.0 with the LES simulation in section 2.1. Furthermore, c was found to be close to 3.0 across all the numerical experiments from section 3.2. Substituting this expression in the buoyancy budget, equation 7, we find that,

$$\frac{\Delta b}{h N^2} \simeq \left(\frac{1}{2} - \frac{1}{c} \right). \quad (25)$$

This expression can then be used as a new boundary layer depth criterion that replaces equation 13,

$$C^* = \frac{h \left(\frac{1}{C^{sh}} \int_{-C^{sh}}^0 B(z) dz - B(-h) \right)}{N^2 h^2 + 10^{-11} \text{m}^2 \text{s}^{-2}}, \quad (26)$$

where C^* replaces C^H as the dimensionless parameter whose value sets the boundary layer depth. Based on equation 25, we expect

$$C^* \simeq \left(\frac{1}{2} - \frac{1}{c} \right) \simeq \frac{1}{6}, \quad (27)$$

321 based on the LES result. The relation equation 26 is an implicit equation for h which
 322 guarantees that equation 24 holds. Once again, it may be useful to point out that C^*
 323 takes the place of what is generally referred to as the bulk Richardson number in the KPP
 324 literature, but that nomenclature is inappropriate for the case of penetrative convection
 325 where C^* parameterizes the effect of convective entrainment rather than shear mixing
 326 at the base of the mixed layer.

327 We now repeat the model calibration in section 3.2 with this new boundary layer
 328 depth criterion to test whether there is an optimal value of C^* that is independent of
 329 initial stratification. We estimate all KPP parameters and show the new mixing depth
 330 parameter for simulations with different initial stratifications in figure 7. There is no ob-
 331 vious trend in the optimal values of C^* and the error bars overlap for all cases. This val-
 332 idates the new criterion in that parameters estimated in different regimes are now con-
 333 sistent with one another. The uncertainties in C^* translate into an uncertainty in bound-
 334 ary layer depth prediction. In particular, values between $0.05 \leq C^* \leq 0.2$ imply a bound-
 335 ary layer depth growth in the range $\sqrt{2.22tQ_b/N^2} \leq h \leq \sqrt{3.33tQ_b/N^2}$.

336 Additionally one can check if the constants estimated with the methodology of sec-
 337 tion 3 are consistent with an *independent* measure directly from the diagnosed LES sim-
 338 ulation. In particular the LES simulations suggest that $C^* \simeq 1/6$ as per equation 27.
 339 From figure 7 we see that the optimal C^* is somewhat smaller than $1/6 = 0.167$ (the
 340 dashed black line). A reason for this discrepancy is the neglect of curvature in the buoy-
 341 ancy budget, since we assumed a piece-wise linear buoyancy profile. Another one is the
 342 finite resolution in the model. A systematic source of error is how we diagnose the bound-
 343 ary layer depth: a different definition, such as the depth of maximum stratification, would
 344 yield a different scaling law (but still proportional to \sqrt{t}). At any rate the Bayesian pa-
 345 rameter estimation bypasses these ambiguities / inconsistencies by direct comparison with
 346 the LES data.

347 We do not explore other modifications to the boundary layer depth criterion as this
 348 would greatly expand the scope of this article. The criterion described in this section as-
 349 sumes a constant initial stratification and a constant surface heat loss, which leads to
 350 the \sqrt{t} growth of the boundary layer depth. It would be interesting to extend the cri-
 351 terion to arbitrary initial stratification, variable surface heat fluxes, not to mention the
 352 interaction with wind-driven mixing. The goal here was not to derive a new parameter-
 353 ization, but rather to introduce a methodology for obtaining meaningful parameteriza-
 354 tions for climate models.

355 4 Discussion

356 In this work we have used a Bayesian methodology for estimating parameters in
 357 parameterizations of subgrid-scale physics as a first step towards parameter sensitivity
 358 studies for Earth Systems Models. We have calculated parameter probability distribu-
 359 tions for parameters in the K -profile parametrization (KPP) by comparing with very high
 360 resolution simulations of ocean convection.

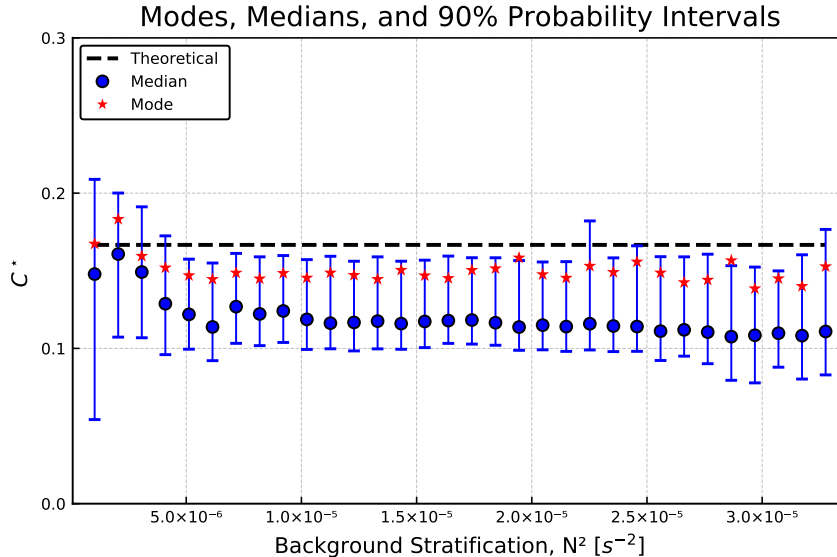


Figure 7. The modified mixing depth parameter optimized across various background stratification. The dots are the median values, the stars are the mode, and the error bars correspond to 90% probability intervals. The dashed line corresponds to $1/6$, the theoretical expectation based on equation 27. This is similar to figure 6, but using the modification from section 3.3. Here one can see that there mixing depth parameter when estimated across various regimes produces similar results. This is a desirable feature in a parameterization.

361 This approach differs from current practice in ocean and climate modelling. Stan-
 362 dard practice is to estimate parameters from a high resolution simulation or field cam-
 363 paign, or to tune parameters to reduce biases in global simulations. For example, the nondi-
 364 mensional amplitude of the KPP eddy diffusivity would be estimated as the ratio of the
 365 flux to the gradient from a single or a few high resolution simulations, (Van Roekel et
 366 al., 2018), or field campaigns, (Large et al., 1994). This assumes at the outset that the
 367 parameters calibrated for one test case will apply to all other scenarios, something that
 368 should be demonstrated rather than assumed. The other common approach is to tune
 369 the parameters in global models to reduce biases in climate relevant variables like ocean
 370 heat uptake or sea surface temperature (Menemenlis et al., 2005; Sraj et al., 2016). This
 371 can result in parameter choices that are inconsistent with the subgrid-scale physics they
 372 are supposed to parameterize. Our approach, instead, relies on a suite of high resolu-
 373 tion simulations that span all the scenarios the parameterization is supposed to capture.
 374 Applying a Bayesian methodology, we then estimate the probability distributions for pa-
 375 rameters which are consistent with the whole suite of high resolution simulations. It is
 376 worth pointing out that the methodology is computationally trivial once one has the LES
 377 solutions. The intellectual effort goes into identifying appropriate forms for the cost func-
 378 tions and probability distributions to guide the quantification of parameter values and
 379 their uncertainty.

380 We illustrated our approach to estimating KPP parameters for convection in a strat-
 381 ified ocean. We found that no unique set of parameters could capture the deepening of
 382 convection for different initial stratifications. We showed that a reformulation of the cri-
 383 terion to estimate the penetration depth of convection allowed us to find parameters that
 384 agreed well with the whole set of high resolution simulations. This shows the Bayesian
 385 approach is not only useful to estimate the probability distributions of parameter val-

386 ues in a parameterization, but it can also be used to identify and eliminate potential bi-
 387 ases in parameterizations.

388 Ultimately, the hope is that parameter probability distributions estimated in lo-
 389 cal regimes will be useful for estimating uncertainties in global climate models; however,
 390 when coupling different components together, new parameters are introduced and non-
 391 linear interactions between parameters can arise. Thus additional optimization is required
 392 for the full system and this requires innovation, because the methodologies described in
 393 this paper are not computationally feasibly when applied to larger systems. A promis-
 394 ing approach for the global climate system is the Calibrate, Emulate, and Sample (CES)
 395 philosophy as outlined in (Cleary et al., 2020). In the CES approach, one uses a prior
 396 distribution for parameters (such as the ones calculated here for KPP) in a climate model
 397 in order to generate a preliminary ensemble of parameters. One then evolves this ensem-
 398 ble according to a loss function (appropriate to the global model) to generate a set of
 399 points that serve as a good “nodes” for interpolation of the loss function (or, alterna-
 400 tively, the “forward map”). A model, also called the emulator/surrogate model, is then
 401 chosen as an interpolator: this can be, for example, a Bayesian Neural Network or a Gaus-
 402 sian Process. Then one uses the interpolated function to calculate the probability of the
 403 posterior distribution using classic algorithms like the RW-MCMC method. In this way,
 404 one avoids rerunning the climate model and instead leverages as much information as
 405 possible from limited data. The surrogate model can then be used to update the prior
 406 distribution and improve predictions of the global model.

407 Stated differently we do not allow for arbitrary parameter perturbations when try-
 408 ing to match a climate model to data, parameter perturbations must take into account
 409 prior information. We propose obtaining this prior information by using highly resolved
 410 local simulations of turbulence. These experiments must be carefully designed and take
 411 into account suites of subgrid scale processes that one might expect to encounter in a
 412 global ocean model: vertical mixing, baroclinic effects, Langmuir turbulence, surface wave
 413 effects, bottom boundary layer turbulence, etc. If the global problem still exhibits sig-
 414 nificant biases after using all available prior information then this suggests that there
 415 is a fundamental deficiency in our understanding of how the different components of the
 416 climate system interact with one another. In this way one can start decoupling where
 417 biases in climate models come from. When the physics of local processes are well under-
 418 stood, then the additional uncertainties induced by coupling different regimes can be iso-
 419 lated and scrutinized. One no longer allows for biases to compensate for one another.

420 Appendix A Oceananigans.jl

Oceananigans.jl is open source software for ocean process studies written in the Ju-
 lia programming language (Bezanson et al., 2017; Ramadhan et al., 2020; Besard et al.,
 2019). For the large eddy simulations (LESs) reported in this paper, Oceananigans.jl is
 configured to solve the spatially-filtered, incompressible Boussinesq equations with a tem-
 perature tracer equations. Letting $\mathbf{u} = (u, v, w)$ be the three-dimensional, spatially-filtered
 velocity field, θ be the conservative temperature, p be the kinematic pressure, f be the
 Coriolis parameter, and $\boldsymbol{\tau}$ and \mathbf{q} be the stress tensor and temperature flux due to sub-
 filter turbulent diffusion, the equations of motion are A1–A3,

$$\partial_t \mathbf{u} + (\mathbf{u} \cdot \nabla) \mathbf{u} + f \hat{z} \times \mathbf{u} + \nabla p = b \hat{z} - \nabla \cdot \boldsymbol{\tau}, \quad (\text{A1})$$

$$\partial_t \theta + \mathbf{u} \cdot \nabla \theta = -\nabla \cdot \mathbf{q}, \quad (\text{A2})$$

$$\nabla \cdot \mathbf{u} = 0. \quad (\text{A3})$$

The buoyancy b appearing in A1 is related to conservative temperature by a linear equa-
 tion of state,

$$b = \alpha g (\theta_0 - \theta), \quad (\text{A4})$$

421 where $\theta_0 = 20^\circ \text{C}$ is a reference temperature, $\alpha = 2 \times 10^{-4} (\text{°C})^{-1}$ is the thermal ex-
 422 pansion coefficient, and $g = 9.81 \text{ ms}^{-1}$ is gravitational acceleration at the Earth’s sur-
 423 face.

424 **A1 Subfilter stress and temperature flux**

The subfilter stress and momentum fluxes are modeled with downgradient closures, such that

$$\tau_{ij} = -2\nu_e \Sigma_{ij} \quad \text{and} \quad \mathbf{q} = -\kappa_e \nabla \theta, \quad (\text{A5})$$

425 where $\Sigma_{ij} \equiv \frac{1}{2} (\partial_i u_j + \partial_j u_i)$ is the strain rate tensor, and ν_e and κ_e are the eddy vis-
 426 cosity and eddy diffusivity of conservative temperature. The eddy viscosity ν_e and eddy
 427 diffusivity κ_e in equation A5 are modeled with the anisotropic minimum dissipation (AMD)
 428 formalism introduced by (Rozema et al., 2015) and (Abkar et al., 2016), refined by (Verstappen,
 429 2018), and validated and described in detail for ocean-relevant scenarios by (Vreugdenhil
 430 & Taylor, 2018). AMD is simple to implement, accurate on anisotropic grids (Vreugdenhil
 431 & Taylor, 2018), and relatively insensitive to resolution (Abkar et al., 2016).

432 **A2 Numerical methods**

433 To solve equations A1–A3 with the subfilter model in equation A5 we use the soft-
 434 ware package ‘`Oceananigans.jl`’ written in the high-level Julia programming language
 435 to run on Graphics Processing Units, also called ‘GPUs’ (Bezanson et al., 2017; Besard
 436 et al., 2019; Besard et al., 2019). `Oceananigans.jl` uses a staggered C-grid finite vol-
 437 ume spatial discretization (Arakawa & Lamb, 1977) with centered second-order differ-
 438 ences to compute the advection and diffusion terms in equation A1 and equation A2, a
 439 pressure projection method to ensure the incompressibility of \mathbf{u} , a fast, Fourier-transform-
 440 based eigenfunction expansion of the discrete second-order Poisson operator to solve the
 441 discrete pressure Poisson equation on a regular grid (Schumann & Sweet, 1988), and second-
 442 order explicit Adams-Bashforth time-stepping. For more information about the staggered
 443 C-grid discretization and second-order Adams-Bashforth time-stepping, see section 3 in
 444 (Marshall et al., 1997) and references therein. The code and documentation are avail-
 445 able for perusal at <https://github.com/climate-machine/Oceananigans.jl>.

446 **Appendix B Plume Model Derivation of the Mixing Layer Depth Cri-** 447 **teria**

We begin by considering the vertical momentum equation for a parcel punching through the transition layer,

$$w' \frac{dw'}{dz} \simeq -(b' - \bar{b}) \quad (\text{B1})$$

where b' is the buoyancy of the parcel, assumed to be equal to the mixed layer value and \bar{b} is the area mean buoyancy profile in the transition layer. This equation holds if the mean buoyancy profile is in hydrostatic balance and the area occupied by sinking plumes is small compared to the total area (Deardorff et al., 1980). Integrating from $z = -h + \Delta h$, where $w' \equiv w_e$, to $z = -h$, where the turbulence and particle descent vanish and hence $w' = 0$, gives,

$$(w_e)^2 \simeq N_e^2 \Delta h^2, \quad (\text{B2})$$

assuming that the background stratification N_e^2 is constant in the entrainment layer. Introducing Δb as the difference between the buoyancy in the mixed layer and that at the base of the transition layer, we have, $\Delta \bar{b} = N_e^2 \Delta h$, and hence,

$$\Delta \bar{b} \propto w^* N_e, \quad (\text{B3})$$

where we assumed that $w_e \propto w^*(-h + \Delta h)$. The criterion for diagnosing the mixing layer depth follows from this relationship; h is defined as the first depth z below the ocean surface where,

$$\frac{\Delta b_p(z)}{w^*(z)N_e(z)} = C^H, \quad (\text{B4})$$

448 for some universal constant C^H .

449 **Appendix C A Primer on on Probability Distributions**

450 In defining a probability distribution, there are a few desirable features that enable us to make a direct connection to a loss function, \mathcal{L} :

- 452 1. In the limit of no uncertainty, a probability distribution should collapse to a delta function centered at optimal parameter values of the loss function.
- 453 2. The most probable value of the distribution should correspond to the optimal parameter in a loss function.
- 454 3. The uncertainty of a parameter value should be determined in terms of its effect on the loss function.

455 There are many probability distributions that satisfy the above criteria, but we choose

$$\rho(\mathbf{C}) \propto \rho^0(\mathbf{C}) \exp(-\mathcal{L}(\mathbf{C})/\mathcal{L}_0), \quad (\text{C1})$$

458 where ρ^0 is a uniform prior distribution, \mathcal{L} is a loss function, and \mathcal{L}_0 is a hyperparameter. In the following subsections we hope to elucidate why we made our particular choices. In section C1 we describe why we chose \mathcal{L}_0 as we did. In section C2 we go into detail on how to sample the probability distribution via the RW-MCMC algorithm as well as intuition for what it is doing.

463 **C1 Explanation for our choice of \mathcal{L}_0**

464 We now describe what the parameter \mathcal{L}_0 means in more detail. The limit $\mathcal{L}_0 \rightarrow 0$ corresponds to no uncertainty. In this limit, the probability distribution collapses to a delta function centered around the global optimal parameters of the loss function. An easy way to see this is to interpret the definition of the probability distribution as a Boltzmann-Gibbs distribution where the loss function corresponds to the energy of the system, and the constant \mathcal{L}_0 is analogous to the temperature of the system, kT where k is the Boltzmann constant, and T is the temperature. In the limit of zero temperature, the system collapses to the lowest energy state, in this case, the minimum of the loss function. The alternative limit $\mathcal{L}_0 \rightarrow \infty$ corresponds to an uncertainty that reduces to the prior distribution, ρ^0 . In this case, information gained from loss function evaluations are uninformative. In analogy with the Boltzmann-Gibbs distribution this corresponds to infinite temperature and every energy state becomes equally likely (hence uninformative).

The maximum of the probability distribution is the mode. This value will be independent of \mathcal{L}_0 if the prior distribution is uniform. Indeed, letting \mathbf{C}^* denote the (global) minimum of the loss function and \mathbf{C} denote any other value, we get

$$\mathcal{L}(\mathbf{C}^*) \leq \mathcal{L}(\mathbf{C}) \Rightarrow \exp(-\mathcal{L}(\mathbf{C})/\mathcal{L}_0) \leq \exp(-\mathcal{L}(\mathbf{C}^*)/\mathcal{L}_0) \Rightarrow \rho(\mathbf{C}) \leq \rho(\mathbf{C}^*). \quad (\text{C2})$$

476 Hence the minimum of the loss function is the most probable value of the probability distribution independent of \mathcal{L}_0 for a uniform prior distribution.

478 As mentioned in section 3, we choose the hyperparameter \mathcal{L}_0 to be the minimum of the loss function \mathcal{L} . Let us take a step back and explain *why* we used this definition.

480 Whatever the choice of \mathcal{L}_0 , we would like the uncertainties of parameters to be indepen-
 481 dent of the units for which we use to evaluate the loss function. Furthermore, we would
 482 like for it to become smaller when there is less model bias and greater when there is more
 483 model bias. In other words, there is more uncertainty when the parameterization does
 484 a poor job of matching the “truth”. This second criteria suggests that \mathcal{L}_0 should be
 485 a monotonic function of the global minimum of the loss function (the bias) $\mathcal{L}(\mathbf{C}^*)$ where
 486 \mathbf{C}^* denotes the optimal parameter values. The first criteria, coupled with the fact that
 487 a perfect model would output a value of 0 for the loss function, yields a number of choices.
 488 We use the simple form $\mathcal{L}_0 \propto \mathcal{L}(\mathbf{C}^*)$, as it is consistent with the previous sentence. Stated
 489 differently, we take \mathcal{L}_0 to be proportional to model bias. It is perhaps more correct to
 490 think of \mathcal{L}_0 as corresponding to *differences* in the loss function. Here we are using the
 491 difference between the “best” parameter choice and a (perhaps) non-existent paramete-
 492 r choice that would correspond to a perfect model whose loss function value is zero.

493 In the case that we have a *perfect* model and *perfect* data this would correspond
 494 to no uncertainty in the parameters and we would go back to having a point estimate
 495 for parameter values. Our choice here naturally assumes that both our data and our loss
 496 function are “perfect”. If we have an idea of how imperfect our data or loss function may
 497 be, then this additional uncertainty should be taken into account in the choice of \mathcal{L}_0 , or
 498 more generally, in the functional form of the probability distribution. We do not con-
 499 sider this additional source of uncertainty here.

500 The choice that we have made for this parameter may be thought of as how far we
 501 are willing to deviate from optimal values while still producing small changes in the loss
 502 function. In other words it is an *e*-folding length defined by the minimum of the loss func-
 503 tion. In this sense, we do not prescribe the uncertainty of a parameter a priori, but in-
 504 stead, implicitly determine it from a choice of how far from optimal parameters that we
 505 are willing to deviate. One could spend a lifetime arguing about the finer details of a choice
 506 for this parameter, but ultimately this would detract from the real goal: to gain an un-
 507 derstanding of what happens when one perturbs parameters away from optimal values.
 508 Any choice of $\mathcal{L}_0 > 0$ would do this. The important part is to make clear how such a
 509 choice is made and why.

510 Admittedly, in practice it is seldom possible to find the true global optimum of \mathcal{L}
 511 and the best one could hope for is some approximate value that is the “best known” op-
 512 timal value $\tilde{\mathbf{C}}$ to get an approximate $\tilde{\mathcal{L}}_0 \equiv \mathcal{L}(\tilde{\mathbf{C}})$. Since $\mathcal{L}_0 = \mathcal{L}(\mathbf{C}^*) \leq \mathcal{L}(\tilde{\mathbf{C}}) = \tilde{\mathcal{L}}_0$,
 513 our uncertainty is a conservative estimate (recall that a smaller \mathcal{L}_0 corresponds to *less*
 514 uncertainty).

515 C2 Random Walk Markov Chain Monte Carlo

516 We use standard methods to sample values from the probability distribution. The
 517 algorithm that we describe here is the Random Walk Markov Chain Monte Carlo Method
 518 (RW-MCMC), first used in (Metropolis et al., 1953). This method is most appropriate
 519 for loss functions where gradient information is either unavailable or prohibitively ex-
 520 pensive to calculate. If it is possible to differentiate the loss function, then other meth-
 521 ods may be more efficient at sampling from the distribution, such as Hamiltonian Monte
 522 Carlo. This alternative method is especially relevant when one wants to estimate a large
 523 number of parameters; however, here we are estimating four parameters, and there is no
 524 need for additional complexity.

525 The RW-MCMC algorithm, as the name suggests, performs a random walk in pa-
 526 rameter space \mathbf{C} . It stays in regions of high probability more often, thereby allowing one
 527 to take the trajectories of the random walk and construct histograms that are directly
 528 related to the empirical distribution of underlying probability distribution function.

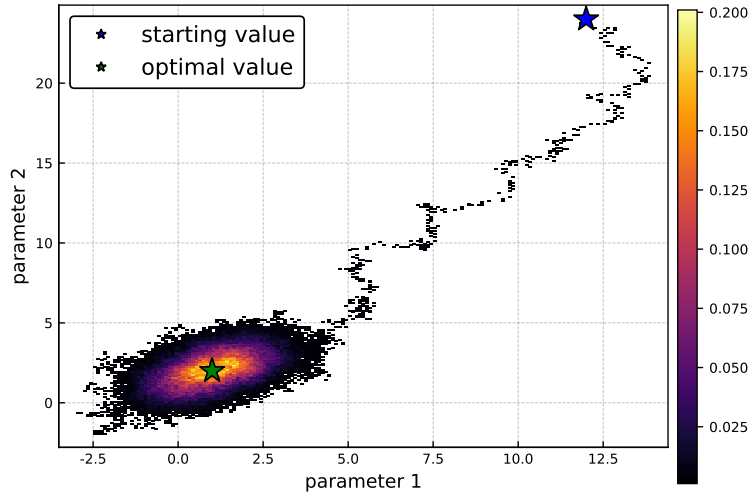


Figure C1. A histogram plot of RW-MCMC output for the target probability distribution equation C3 starting from a suboptimal value $\mathbf{C}^0 = [12.0 \ 24.0]^T$ and using 10^5 iterations. The white space signifies regions that the algorithm did not explore. The dark regions correspond to places that were rarely visited whereas the red and yellow regions correspond to places that were visited more often. The trail from the starting value to the optimal value illustrates how the random walk is biased towards regions of ever increasing probability.

We further illustrate what the algorithm does by considering the following toy loss function:

$$\mathcal{L}(\mathbf{C}) = \frac{1}{2}(\mathbf{C} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{C} - \boldsymbol{\mu}), \quad \boldsymbol{\mu} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \quad \text{and} \quad \Sigma = \begin{bmatrix} 1 & 1/2 \\ 1/2 & 1 \end{bmatrix}, \quad (\text{C3})$$

529 with the choice $\mathcal{L}_0 = 1$ in the probability distribution. In figure C1 we show a typical
530 output from an RW-MCMC algorithm for this probability distribution. Note that the
531 optimal value of the loss function is $\boldsymbol{\mu} = [1 \ 2]^T$. Starting from a poor initial guess, here
532 $\mathbf{C}^0 = [12.0 \ 24.0]^T$, the algorithm goes towards regions of higher probability (lower cost
533 function) by randomly choosing which direction to go. Once a region of high probabili-
534 ty is found, in this case parameter values close to $\boldsymbol{\mu}$, the parameters hover around the
535 minima of the loss function in a way that is consistent with the target probability dis-
536 tribution.

537 We now go into detail on RW-MCMC algorithm. Ratios of the probability distri-
538 bution play a prominent role in the RW-MCMC algorithm; however, due to finite arith-
539 metic considerations it is actually better to work with the logarithms of the probabili-
540 ty distribution. A convenient way to do this is to use the negative log-likelihood func-
541 tion as $\ell = -\ln \rho$. This function can be thought of as being essentially the same as the
542 loss function 15, but shifted and scaled. Denote elements of a sequence of parameter val-
543 ues by \mathbf{C}_i . The RW-MCMC algorithm is:

- 544 1. Choose initial parameter values \mathbf{C}_0 . One choice for this parameter is the best known
545 minimizer of the log-likelihood function using standard minimization techniques.
- 546 2. Calculate a ‘‘proposal parameter’’ $\tilde{\mathbf{C}}_1$. This will be described in more detail later.
- 547 3. Calculate $\Delta\ell = \ell(\mathbf{C}_0) - \ell(\tilde{\mathbf{C}}_1)$. This is a measure of how much more likely $\tilde{\mathbf{C}}_1$
548 is relative to \mathbf{C}_0 .

- 549 4. Draw a random uniform random variable from the interval $[0, 1]$, e.g, calculate $u =$
550 $\mathcal{U}(0, 1)$. This is used to determine whether or not to accept $\tilde{\mathbf{C}}_1$ as a new param-
551 eter.
- 552 5. If $\log(u) < \Delta\ell$ set $\mathbf{C}_1 = \tilde{\mathbf{C}}_1$. Otherwise set $\mathbf{C}_1 = \mathbf{C}_0$. This is the “accept / re-
553 ject” step. Note that if $\Delta\ell > 0$, i.e. the proposed parameter produces a smaller
554 output in the negative log-likelihood function, the proposal is always accepted.
- 555 6. Repeat steps 2-5 for \mathbf{C}_i , replacing $\mathbf{C}_0 \rightarrow \mathbf{C}_i$ and $\mathbf{C}_1 \rightarrow \mathbf{C}_{i+1}$, to generate a se-
556 quence (or chain) of parameter values.

557 Interpreting the negative log-likelihood function as a potential function, the algorithm
558 may be succinctly stated as “always go downhill, sometimes go uphill”. The sequence
559 of parameter values generated by this algorithm can then be used to construct any statis-
560 tics of the probability distribution 16, including empirical distributions, marginal dis-
561 tributions, and joint distributions. In the context of KPP this can also generate the un-
562 certainty of a temperature at a given point in space and time as well as the uncertainty
563 of the mixed layer depth at a given time.

This random walk is different from a random number generator in that successive samples are not independent of one another but are instead correlated. The random walk must be run for enough time to generate a sufficient number of statistically independent samples. The proposal step is crucial to ensure this feature. Thus, we will now describe how to choose a proposal in more detail. If there are no restrictions on the range of parameter values, then one can perturb each parameter by a Gaussian random variable with mean zero and covariance matrix Σ , i.e.

$$\tilde{\mathbf{C}}_{i+1} = \mathbf{C}_i + \mathcal{N}(0, \Sigma) \quad (\text{C4})$$

564 Interestingly, the algorithm is guaranteed to work *independent* of the choice of Σ as long
565 as the covariance matrix Σ is nonzero and the same proposal is used throughout the ran-
566 dom walk⁶; however, suitable choices of Σ can speed up convergence to the probability
567 distribution. At the end of an RW-MCMC run one can diagnose the “number of inde-
568 pendent samples” by using approximations of the correlation length, see Sokal (1997).
569 If Σ is too small then the acceptance rate will be too large since each proposal param-
570 eter is barely any different from the original parameter. Too large of a proposal often
571 yields too low acceptance rates since it is typically easier to propose a parameter asso-
572 ciated with a region of low probability than high probability (thereby making it likely
573 that one is choosing a point that is “uphill” more often than “downhill”). One option
574 is to take Σ to be a diagonal matrix whose diagonal elements are proportional to the square
575 10% of the default parameter values, i.e., the standard deviation of the proposal of each
576 was about 10% and that each parameter component proposal was independent. A com-
577 mon option is to choose Σ according to the covariance matrix of the prior distribution.
578 Yet, another option is to perform a preliminary random walk to estimate the covariance
579 of the target distribution and then use this estimated covariance matrix in a new ran-
580 dom walk. In general, there is no rule that will always speed up convergence, but we found
581 that the last method to gave the best results.

582 If we would like to restrict the parameters to be in a finite range it is as simple as
583 making the random walk take place in a periodic domain in parameter space. Another
584 option is to redefine the loss function so that it outputs infinity if one plugs in a value
585 outside a specified range. We opt for the former since it does not “waste” function eval-
586 uations.

⁶ If one decides to change the proposal one needs to start the random walk over and cannot reuse data generated from another proposal.

Acknowledgments

The authors would like to thank Carl Wunsch, Tapio Schneider, Andrew Stuart, and William Large for numerous illuminating discussions. Our work is supported by the generosity of Eric and Wendy Schmidt by recommendation of the Schmidt Futures program, and by the National Science Foundation under grant AGS-6939393.

References

- Abkar, M., Bae, H. J., & Moin, P. (2016). Minimum-dissipation scalar transport model for large-eddy simulation of turbulent flows. *Physical Review Fluids*, *1*(4), 041701. doi: 10.1103/PhysRevFluids.1.041701
- Albers, D. J., Blancquart, P.-A., Levine, M. E., Seylabi, E. E., & Stuart, A. (2019). Ensemble kalman methods with constraints. *Inverse Problems*, *35*(9), 095007. doi: 10.1088/1361-6420/ab1c09
- Arakawa, A., & Lamb, V. R. (1977). Computational design of the basic dynamical processes of the ucla general circulation model. In J. Chang (Ed.), *General circulation models of the atmosphere* (Vol. 17, p. 173 - 265). Elsevier. doi: 10.1016/B978-0-12-460817-7.50009-4
- Besard, T., Churavy, V., Edelman, A., & Sutter, B. D. (2019). Rapid software prototyping for heterogeneous and distributed platforms. *Advances in Engineering Software*, *132*, 29 - 46. Retrieved from <http://www.sciencedirect.com/science/article/pii/S0965997818310123> doi: <https://doi.org/10.1016/j.advengsoft.2019.02.002>
- Besard, T., Foket, C., & De Sutter, B. (2019, April). Effective extensible programming: Unleashing julia on gpus. *IEEE Transactions on Parallel and Distributed Systems*, *30*(4), 827-841. doi: 10.1109/TPDS.2018.2872064
- Bezanson, J., Edelman, A., Karpinski, S., & Shah, V. B. (2017). Julia: A Fresh Approach to Numerical Computing. *SIAM Review*, *59*(1), 65–98. doi: 10/f9wkpj
- Cleary, E., Garbuno-Inigo, A., Lan, S., Schneider, T., & Stuart, A. M. (2020). Calibrate, emulate, sample. *arXiv:2001.03689 [stat.CO]*. Retrieved from <https://arxiv.org/abs/2001.03689>
- Deardorff, J. W., Willis, G. E., & Stockton, B. H. (1980). Laboratory studies of the entrainment zone of a convectively mixed layer. *Journal of Fluid Mechanics*, *100*(1), 41–64. doi: 10.1017/S0022112080001000
- DuVivier, A. K., Large, W. G., & Small, R. J. (2018). Argo observations of the deep mixing band in the southern ocean: A salinity modeling challenge. *Journal of Geophysical Research: Oceans*, *123*(10), 7599-7617. doi: 10.1029/2018JC014275
- Golaz, J.-C., Larson, V. E., Hansen, J. A., Schanen, D. P., & Griffin, B. M. (2007). Elucidating model inadequacies in a cloud parameterization by use of an ensemble-based calibration framework. *Monthly Weather Review*, *135*(12), 4077-4096. doi: 10.1175/2007MWR2008.1
- Hourdin, F., Mauritsen, T., Gettelman, A., Golaz, J.-C., Balaji, V., Duan, Q., ... Williamson, D. (2017). The art and science of climate model tuning. *Bulletin of the American Meteorological Society*, *98*(3), 589-602. doi: 10.1175/BAMS-D-15-00135.1
- Kara, A. B., Rochford, P. A., & Hurlburt, H. E. (2000). An optimal definition for ocean mixed layer depth. *Journal of Geophysical Research: Oceans*, *105*(C7), 16803-16821. doi: 10.1029/2000JC900072
- Kirkpatrick, S., Gelatt, C. D., & Vecchi, M. P. (1983). Optimization by simulated annealing. *Science*, *220*(4598), 671–680. doi: 10.1126/science.220.4598.671
- Large, W. G., McWilliams, J. C., & Doney, S. C. (1994). Oceanic vertical mixing: A review and a model with a nonlocal boundary layer parameterization. *Reviews of Geophysics*, *32*(4), 363-403. doi: 10.1029/94RG01872

- 640 Large, W. G., Patton, E. G., DuVivier, A. K., Sullivan, P. P., & Romero, L. (2019).
641 Similarity theory in the surface layer of large-eddy simulations of the wind-,
642 wave-, and buoyancy-forced southern ocean. *Journal of Physical Oceanography*,
643 *49*(8), 2165-2187. doi: 10.1175/JPO-D-18-0066.1
- 644 Marshall, J., Adcroft, A., Hill, C., Perelman, L., & Heisey, C. (1997). A finite-
645 volume, incompressible navier stokes model for studies of the ocean on parallel
646 computers. *Journal of Geophysical Research: Oceans*, *102*(C3), 5753-5766. doi:
647 10.1029/96JC02775
- 648 Marshall, J., & Schott, F. (1999). Open-ocean convection: Observations, theory, and
649 models. *Reviews of Geophysics*, *37*(1), 1-64. doi: 10.1029/98RG02739
- 650 Menemenlis, D., Fukumori, I., & Lee, T. (2005, 05). Using green's functions to cali-
651 brate an ocean general circulation model. *Monthly Weather Review*, *133*, 1224-
652 1240. doi: 10.1175/MWR2912.1
- 653 Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., & Teller, E.
654 (1953). Equation of state calculations by fast computing machines. *The*
655 *Journal of Chemical Physics*, *21*(6), 1087-1092. doi: 10.1063/1.1699114
- 656 Ramadhan, A., Wagner, G. L., Hill, C., Campin, J.-M., Churavy, V., Besard,
657 T., ... Ferrari, R. (2020). Oceananigans.jl: Fast and friendly geo-
658 physical fluid dynamics on GPUs. *The Journal of Open Source Soft-*
659 *ware*, *4*(44), 2018. Retrieved from [https://joss.theoj.org/papers/
660 aeb75d5b99d40ae6e0c8a3a4f09f3285](https://joss.theoj.org/papers/aeb75d5b99d40ae6e0c8a3a4f09f3285) (Under review.)
- 661 Rozema, W., Bae, H. J., Moin, P., & Verstappen, R. (2015). Minimum-dissipation
662 models for large-eddy simulation. *Physics of Fluids*, *27*(8), 085107. doi: 10
663 .1063/1.4928700
- 664 Schneider, T., Lan, S., Stuart, A., & Teixeira, J. (2017). Earth system modeling
665 2.0: A blueprint for models that learn from observations and targeted high-
666 resolution simulations. *Geophysical Research Letters*, *44*(24), 12,396-12,417.
667 doi: 10.1002/2017GL076101
- 668 Schneider, T., Teixeira, J., Bretherton, C., Brient, F., Pressel, K., Schär, C., &
669 Siebesma, A. (2017, 01). Climate goals and computing the future of clouds.
670 *Nature Climate Change*, *7*, 3-5. doi: 10.1038/nclimate3190
- 671 Schumann, U., & Sweet, R. A. (1988). Fast fourier transforms for direct solution
672 of poisson's equation with staggered boundary conditions. *Journal of Compu-*
673 *tational Physics*, *75*(1), 123 - 137. doi: [https://doi.org/10.1016/0021-9991\(88\)
674 90102-7](https://doi.org/10.1016/0021-9991(88)90102-7)
- 675 Sokal, A. (1997). Monte carlo methods in statistical mechanics: Foundations
676 and new algorithms. In C. DeWitt-Morette, P. Cartier, & A. Folacci (Eds.),
677 *Functional integration: Basics and applications* (pp. 131-192). Boston, MA:
678 Springer US. doi: 10.1007/978-1-4899-0319-8_6
- 679 Sraj, I., Zedler, S. E., Knio, O. M., Jackson, C. S., & Hoteit, I. (2016). Polyno-
680 mial chaos-based bayesian inference of k-profile parameterization in a general
681 circulation model of the tropical pacific. *Monthly Weather Review*, *144*(12),
682 4621-4640. Retrieved from <https://doi.org/10.1175/MWR-D-15-0394.1> doi:
683 10.1175/MWR-D-15-0394.1
- 684 Van Roekel, L., Adcroft, A. J., Danabasoglu, G., Griffies, S. M., Kauffman, B.,
685 Large, W., ... Schmidt, M. (2018). The kpp boundary layer scheme for the
686 ocean: Revisiting its formulation and benchmarking one-dimensional simula-
687 tions relative to les. *Journal of Advances in Modeling Earth Systems*, *10*(11),
688 2647-2685. doi: 10.1029/2018MS001336
- 689 Verstappen, R. (2018). How much eddy dissipation is needed to counter-
690 balance the nonlinear production of small, unresolved scales in a large-
691 eddy simulation of turbulence? *Computers & Fluids*, *176*, 276-284. doi:
692 10.1016/j.compfluid.2016.12.016
- 693 Vreugdenhil, C. A., & Taylor, J. R. (2018). Large-eddy simulations of stratified
694 plane Couette flow using the anisotropic minimum-dissipation model. *Physics*

