



RESEARCH ARTICLE

10.1029/2024MS004485

Online Learning of Entrainment Closures in a Hybrid Machine Learning Parameterization

Costa Christopoulos¹ , Ignacio Lopez-Gomez^{1,2} , Tom Beucler^{3,4} , Yair Cohen^{1,5} , Charles Kawczynski¹, Oliver R. A. Dunbar¹ , and Tapio Schneider¹ 

¹California Institute of Technology, Pasadena, CA, USA, ²Now at Google Research, Mountain View, CA, USA, ³Faculty of Geosciences and Environment, University of Lausanne, Lausanne, Switzerland, ⁴Expertise Center for Climate Extremes, University of Lausanne, Lausanne, Switzerland, ⁵Now at NVIDIA Corporation, Santa Clara, CA, USA

Key Points:

- We train a hybrid subgrid parameterization to minimize the mismatch between a single-column model and large-eddy simulation mean states
- Within the parameterization, the entrainment mixing closure is fully data-driven and trained online via ensemble Kalman inversion
- With no prior information on entrainment, we learn physically realistic mixing closures indirectly from mean simulation states

Correspondence to:

C. Christopoulos,
cchristo@caltech.edu

Citation:

Christopoulos, C., Lopez-Gomez, I., Beucler, T., Cohen, Y., Kawczynski, C., Dunbar, O. R. A., & Schneider, T. (2024). Online learning of entrainment closures in a hybrid machine learning parameterization. *Journal of Advances in Modeling Earth Systems*, 16, e2024MS004485. <https://doi.org/10.1029/2024MS004485>

Received 31 MAY 2024

Accepted 21 OCT 2024

Abstract This work integrates machine learning into an atmospheric parameterization to target uncertain mixing processes while maintaining interpretable, predictive, and well-established physical equations. We adopt an eddy-diffusivity mass-flux (EDMF) parameterization for the unified modeling of various convective and turbulent regimes. To avoid drift and instability that plague offline-trained machine learning parameterizations that are subsequently coupled with climate models, we frame learning as an inverse problem: Data-driven models are embedded within the EDMF parameterization and trained online in a one-dimensional vertical global climate model (GCM) column. Training is performed against output from large-eddy simulations (LES) forced with GCM-simulated large-scale conditions in the Pacific. Rather than optimizing subgrid-scale tendencies, our framework directly targets climate variables of interest, such as the vertical profiles of entropy and liquid water path. Specifically, we use ensemble Kalman inversion to simultaneously calibrate both the EDMF parameters and the parameters governing data-driven lateral mixing rates. The calibrated parameterization outperforms existing EDMF schemes, particularly in tropical and subtropical locations of the present climate, and maintains high fidelity in simulating shallow cumulus and stratocumulus regimes under increased sea surface temperatures from AMIP4K experiments. The results showcase the advantage of physically constraining data-driven models and directly targeting relevant variables through online learning to build robust and stable machine learning parameterizations.

Plain Language Summary In this research, we aim to improve projections of the Earth's climate response by creating a hybrid model that integrates machine learning (ML) into parts of an existing atmospheric model that are less certain. This integration improves our hybrid model's performance, particularly in tropical and subtropical oceanic regions. Unlike previous approaches that first trained the ML and then ran the host model with ML embedded, we train the ML while the host model is running in a single column, which makes the model more stable and reliable. Indeed, when tested under conditions with higher sea surface temperatures, our model accurately predicts outcomes even in scenarios that were not encountered during the ML training. Our study highlights the value of combining ML and traditional atmospheric models for more robust and data-driven climate predictions.

1. Introduction

The latest suite of global climate models (GCMs) continues to exhibit a large range of climate sensitivities, the measure of Earth's equilibrium temperature response to a doubling of atmospheric greenhouse gas concentrations (Meehl et al., 2020). Variance in modeled responses has been traced to disparate representations of subgrid-scale (SGS) processes not explicitly resolved by climate models, specifically those controlling the characteristics of cloud feedbacks (Bony et al., 2015; Sherwood et al., 2014; Vial et al., 2013; Zelinka et al., 2020). Furthermore, climate models often fail to reproduce several key statistics from the recent past when run retrospectively (Vignesh et al., 2020). In light of these discrepancies, researchers have launched systematic efforts across the climate modeling enterprise to incorporate machine learning (ML) methods into GCMs, in order to improve the ability of climate model components to learn from high fidelity data. This study specifically uses a training data set focused on marine low cloud regimes in the central and eastern Pacific—areas that are particularly problematic to model in GCMs (Nam et al., 2012; Črnivec et al., 2023), yet are critical for precise assessments of equilibrium climate sensitivity due to cloud feedbacks (Brient & Schneider, 2016; Myers et al., 2021; Siler et al., 2018).

© 2024 The Author(s). Journal of Advances in Modeling Earth Systems published by Wiley Periodicals LLC on behalf of American Geophysical Union. This is an open access article under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

Initiatives to replace existing physics-based parameterizations in atmospheric models entirely with ML are often marred with challenges surrounding numerical instability and extrapolation performance. Instabilities, such as the generation of unstable gravity wave modes (Brenowitz et al., 2020), largely arise from feedbacks between the learned SGS parameterization and the dynamical core upon integration. Currently, the favored strategy is to train ML models offline via supervised learning to predict SGS tendencies as a function of the resolved atmospheric state, then couple trained models to a dynamical core to perform inferences at each model timestep (Krasnopolsky et al., 2013; Rasp et al., 2018; Yuval & O’Gorman, 2020). As an example of the offline training procedure for atmospheric turbulence, a recent encoder-decoder approach was used to learn vertical turbulent fluxes in dry convective boundary layers on the basis of coarse-grained large-eddy simulations (Shamekh & Gentine, 2023). Although significant progress has been made toward advancing and stabilizing data-driven parameterizations (Brenowitz & Bretherton, 2019; Wang et al., 2022; Watt-Meyer et al., 2023), the conventional offline training strategy precludes learning unobservable processes indirectly from relevant climate statistics. Furthermore, instabilities arising from system feedbacks are not typically incorporated into training, and cannot be easily assessed until ML models are coupled to a dynamical core (Ott et al., 2020; Rasp, 2020). More recently, the advent of differentiable simplified general circulation models (e.g., without phase transitions of water) has enabled spatially three-dimensional (3D) online training of ML-based SGS parameterizations using short-term forecasts (Kochkov et al., 2024). These strategies have not yet overcome the problems of instability and extrapolation to warmer climates and remain difficult to interpret.

We take steps to address these issues by employing ensemble Kalman inversion (EKI) to perform parameter estimation within a SGS parameterization from statistics of atmospheric profiles in a single column setup (Dunbar et al., 2021; Huang, Schneider, & Stuart, 2022; M. A. Iglesias et al., 2013). Treating learning as an inverse problem directly enables online learning. Inverse problems are characterized by setups where the dependent variable of some target process is neither directly observable nor explicitly included in the loss function. In this case, it is through secondary causal effects of atmospheric dynamics on observable atmospheric quantities that parameters are optimized. In the field of dynamical systems, theory underpinning the use of inversion techniques to infer parameters is well established (Huang, Huang, et al., 2022; M. A. Iglesias et al., 2013), and they have also been shown to be effective for learning neural networks (NNs), especially in chaotic system where the smoothing properties of ensemble methods can be advantageous (Dunbar et al., 2022; Kovachki & Stuart, 2019). In practice, ensemble Kalman methods have been used to learn drift and diffusion terms in the Lorenz ’96 model (Schneider et al., 2021), nonlinear eddy viscosity models for turbulence (Zhang et al., 2022), the effects of truncated variables in a quasi-geostrophic ocean-atmosphere model (Brajard et al., 2021), and NN-based parameterizations of the quasi-biennial oscillation and gravity waves (Pahlavan et al., 2024). An alternative approach to online learning relies on differentiable methods to explicitly compute gradients through the physical model to learn data-driven components (C. Shen et al., 2023; Um et al., 2021). The differentiable learning approach has been used successfully to learn NN-based closures in numerous idealized turbulence setups (Kochkov et al., 2021; List et al., 2022; MacArt et al., 2021; Shankar et al., 2023). In an Earth system modeling setting, differentiable online learning has been used to learn stable turbulence parameterizations in an idealized quasi-geostrophic setup (Frezat et al., 2022) and residual corrections to an upper-ocean convective adjustment scheme (Ramadhan et al., 2023). While promising, differentiable methods preclude computing gradients through physical models with non-differentiable components, such as the physics stemming from water phase changes in cloud parameterizations. Furthermore, given existing work surrounding differentiable and inverse methods for geophysical fluid dynamics, there remains a lack of literature demonstrating indirect learning of data-driven components in more comprehensive atmospheric parameterizations of convection, turbulence, and clouds. Our contribution is the application of these methods in a more realistic climate modeling setting, a use case which can directly improve operational Earth system models.

We extend a flexible and modular framework that allows for the selective addition of expressive, non-parametric components where physical knowledge is limited, introduced by Lopez-Gomez et al. (2022). Our approach promotes generalizability and interpretability. Interpretability comes by virtue of targeting specific physical processes, which enables a mechanistic analysis of their effect on climate. Generalizability is a result of both retaining this physical framework and employing an inversion strategy that targets climate statistics. The physical framework includes the partial differential equations in which the closure is embedded, the non-dimensionalization of data-driven input variables, and the dimensional scales that modulate learned nondimensional closures. In contrast, a fully data-driven parameterization benefits from expressivity at the expense of

Online Function Learning with Ensemble Kalman Inversion

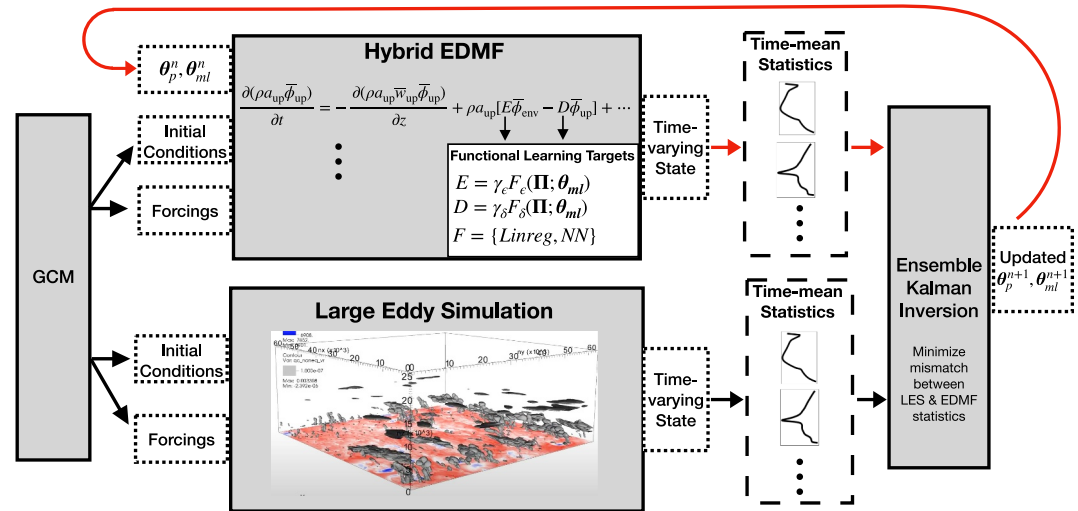


Figure 1. Schematic illustrating the ensemble Kalman inversion pipeline used for online training of a one-dimensional (1D) atmospheric model with both physics-based and data-driven components (hybrid EDMF). Black arrows indicate fixed operations between components, and red arrows indicate dynamic information flow on the basis of Kalman updates to EDMF parameters. The training data comprises 176 LES simulations from the AMIP climate, processed in batches of 16 cases for each ensemble Kalman iteration. Lateral mixing rates are formulated as the product of a dimensional scale γ and a data-driven, nondimensional function F .

sensitivity to training data, leading to difficulties in extrapolating to unobserved climates. Generalizability is verified in our setup by assessing performance on an out-of-distribution climate where SSTs are uniformly increased by 4 K; test error decreases in lockstep with training error from the present climate and overfitting is not observed.

In this study, we will investigate the performance of a single column model containing data-driven lateral mixing closures spanning a range of complexities, from linear regression models to neural networks. In Section 2, we describe in detail the data-driven architectures, training data, and online calibration pipeline. Section 3 outlines the performance of the data-driven eddy-diffusivity mass-flux (EDMF) scheme in terms of the root mean squared error of the mean atmospheric state in a current and warmer climate, and representative vertical profiles are presented with physical implications discussed. Relative to the previous work of Lopez-Gomez et al. (2022), modeling improvements are made by both modifying the calibration pipeline and addressing structural biases in the EDMF model itself, namely boundary conditions and the lateral mixing formulation.

2. Online Training Setup

An overarching goal of SGS modeling is to produce computationally efficient schemes that emulate expensive high-resolution simulations, given the same large-scale forcings, boundary conditions, and initial conditions. Of primary importance are the prediction of SGS fluxes and cloud properties, which are determined by small-scale processes not resolvable by the GCM dynamical core. In the setup described here, parameters in a full-complexity SGS scheme are systematically optimized through the ensemble Kalman inversion technique to match characteristics of high-resolution simulations, namely time-mean vertical profiles and vertically integrated liquid water content produced by large-eddy simulations (LES) (Z. Shen et al., 2022). A variant of the SGS scheme is introduced, which imposes fewer assumptions and incorporates more general data-driven functions that can be determined with data. The SGS model is an eddy-diffusivity mass-flux (EDMF) scheme that parameterizes the effects of turbulence, convection, and clouds. The reference high-resolution simulations are performed with PyCLES (Pressel et al., 2015), which explicitly models convection and turbulent eddies larger than $O(10 \text{ m})$. The process diagram in Figure 1 illustrates how calibrations are performed using the SGS model. Components of the diagram are detailed in the sections that follow, starting with the EDMF scheme.

2.1. Eddy-Diffusivity Mass-Flux (EDMF) Scheme Overview

EDMF schemes partition GCM grid boxes into two or more subdomains, each characterized by containing either coherent structures (updrafts) or relatively isotropic turbulence (environment). While most SGS schemes use separate parameterizations for the boundary layer, shallow convection, deep convection, and stratocumulus regimes, the extended EDMF scheme we use (herein referred to as EDMF) simulates all regimes in a unified manner by making fewer simplifying assumptions (Thuburn et al., 2018). The scheme includes partial differential equations (PDEs) for prognostic updraft properties (notably temperature, humidity, area fraction, and mass flux), which are coupled to PDEs for environmental variables (temperature, humidity, and turbulent kinetic energy). The physical skeleton of the EDMF consists of these coarse-grained equations of motion and houses a collection of closures, appearing as right-hand-side tendency terms for the prognostic variable equations.

The EDMF scheme we use was initially introduced by Tan et al. (2018). It contains closure functions, for example, for entrainment and detrainment, which capture physics without a known, closed-form expression; specifying them is necessary to fully define the set of EDMF PDEs such that they can be numerically integrated. Closures in the EDMF equations play a role similar to SGS parameterizations in grid-scale prognostic equations. Tendencies from SGS parameterizations appear in dynamical core equations, and, similarly, tendencies from closures appear in the EDMF equations. In the context of GCMs, the EDMF parameterization predicts vertical SGS fluxes and cloud properties due to unresolved processes. The present EDMF parameterization, which is run at 50 m vertical resolution, has been shown to effectively generalize between isotropic and stretched vertical grids (Lopez-Gomez et al., 2022). Its prediction of second-order quantities such as turbulent kinetic energy (TKE), which approach zero as the resolution increases, and its inherent SGS memory endow it with some “scale-aware” properties that become especially important as convection begins to be partially resolved in the “gray zone” (Boutle et al., 2014; Schneider et al., 2024; Tan et al., 2018); however, we have not explicitly tested its resolution dependence yet.

Following domain decomposition, the contributions of EDMF SGS fluxes ($\rho \langle w' \phi' \rangle_{\text{sgs}}$) to the grid-scale equation for a general quantity ϕ are

$$\rho \langle w' \phi' \rangle_{\text{sgs}} = -\rho a_{\text{env}} K_{\phi, \text{env}} \frac{\partial \bar{\phi}_{\text{env}}}{\partial z} + \rho a_{\text{up}} (\bar{w}_{\text{up}} - \langle w \rangle) (\bar{\phi}_{\text{up}} - \langle \phi \rangle). \quad (1)$$

Here, $\langle \cdot \rangle$ indicates a grid-mean quantity and $\bar{(\cdot)}$ a subdomain mean. Subscripts “up” and “env” signify updraft and environmental subdomain properties, respectively. We define ρ as the air density, a as the subdomain area fraction, $K_{\phi, \text{env}}$ as the environmental diffusivity of quantity ϕ , and w as the vertical velocity. The first term parameterizes the turbulent flux due to eddy diffusion (ED) in the environment; the second term represents the mass flux (MF) from coherent updrafts. The environmental eddy diffusivity, which governs the diffusive flux, is determined by a mixing length closure (Lopez-Gomez et al., 2020) and environmental TKE, following Mellor and Yamada (1982). In shallow maritime regimes, the turbulent kinetic energy budget is dominated by a balance between shear, buoyancy, and viscous dissipation (Heinze et al., 2015). Thus, lateral mixing primarily affects the updraft mass flux term.

2.1.1. Baseline EDMF: EDMF-20

We compare a hybrid EDMF, detailed in the next section, to a baseline version we call the EDMF-20. The EDMF-20 model includes physically motivated closures for eddy diffusivity (Lopez-Gomez et al., 2020), entrainment/detrainment (Cohen et al., 2020), and perturbation pressure. The physically motivated closure functions were manually tuned so that the simulated EDMF profiles closely match field campaigns. Parameters in EDMF-20 were tuned to match field campaigns representing a spectrum of convective and turbulent regimes, including Bomex (marine shallow convection) (Holland & Rasmusson, 1973), TRMM (deep convection) (Grabowski et al., 2006), a dry convective boundary layer (Soares et al., 2004), ARM-SGP (continental shallow convection) (Brown et al., 2002), RICO (precipitating shallow cumulus) (vanZanten et al., 2011), and DYCOMS (drizzling stratocumulus) (Ackerman et al., 2009; Stevens et al., 2003).

2.1.2. Hybrid EDMF

Building on the baseline EDMF-20, two notable modifications have been implemented since to improve the realism and relax assumptions imposed by previous bottom boundary specifications. First, the surface Dirichlet boundary condition on area fraction, a free parameter found in previous work (Lopez-Gomez et al., 2022) to be correlated with numerous other EDMF parameters, is modified to be a free boundary condition (Appendix A1 in Appendix A). The modification allows updrafts to be generated directly by entrainment and detrainment source terms, rather than being “pinned” to the surface, and eliminates the dependence on lower boundary specification of mass flux and area fraction required by most mass-flux schemes. Second, the surface Dirichlet boundary condition on TKE in previous versions is replaced by a TKE flux boundary condition that depends on surface conditions and turbulence parameters (Appendix A2 in Appendix A).

The key distinction between the hybrid EDMF and EDMF-20 lies in the formulation of data-driven entrainment closures. We consider an EDMF scheme that uses linear regression to determine entrainment rates, designated EDMF-Linreg, and an EDMF scheme that uses a neural network for entrainment rates, designated EDMF-NN. These data-driven closures take the place of the semi-empirical but physically motivated closures implemented in EDMF-20 (Cohen et al., 2020).

2.2. Functional Learning for Entrainment and Detrainment

2.2.1. Functional Learning Targets

Entrainment and detrainment are two forms of cloud mixing, which describe the exchange of mass, momentum, and tracers between coherent updrafts and their turbulent environment (de Rooy et al., 2013). Entrainment is the process whereby environmental properties are incorporated into updrafts, whereas detrainment describes the ejection of updraft properties into the environment. Entrainment and detrainment appear as rates (units of s^{-1}) in the EDMF tendency equations. These processes are often decomposed into the sum of turbulent and dynamical contributions, which represent cloud mixing driven by horizontal turbulent mixing from eddies and exchange due to more organized cloud-scale flows, respectively (de Rooy & Pier Siebesma, 2010). The closures learned for this study combine the contributions into a single function. Inputs for data-driven closures are chosen to be nondimensional variables $\mathbf{\Pi}$. For the closure formulation, we adopt the approach of learning a nondimensional function, which modulates a dimensional scale of the same units as the entrainment/detrainment rates:

$$E = \gamma_e F_e(\mathbf{\Pi}; \Theta_{ml}), \quad (2a)$$

$$D = \gamma_\delta F_\delta(\mathbf{\Pi}; \Theta_{ml}). \quad (2b)$$

Here, γ_e and γ_δ are inverse time scales while F_e and F_δ are nondimensional functions for entrainment and detrainment, respectively. The data-driven functions F parameterize the relationship between nondimensional groups $\mathbf{\Pi}$ and nondimensional mixing rates, given a vector of learnable parameters Θ_{ml} . We note that F_e and F_δ are vertically local functions, and thus map $\mathbf{\Pi}$ groups defined from local quantities at some level to a single lateral mixing rate at that level. Thus, applying the local function to every level yields a vertical profile of mixing rates that varies with height.

The entrainment dimensional scale is chosen as the ratio of updraft-environment vertical velocity difference $\Delta\bar{w}$ to height z :

$$\gamma_e(z) = \frac{\Delta\bar{w}}{z}. \quad (3a)$$

We denote the difference between subdomains with the symbol Δ . Thus, the difference between the mean updraft and environmental vertical velocity is $\Delta\bar{w} = \bar{w}_{up} - \bar{w}_{env}$. The inverse height scaling is chosen here as an easy-to-diagnose proxy of the inverse updraft radius or eddy size at a given height (Siebesma et al., 2007). Thus, γ_e defines a horizontal shear that gives rise to entrainment (Griewank et al., 2022). For detrainment, γ_δ is chosen as a dimensional scale that corresponds to the rate needed to sustain mass flux profiles in steady-state. Taking the EDMF continuity equation (Equation A1) as steady and assuming no horizontal convergence or entrainment yields the detrainment expression

$$\gamma_{\delta}(z) = \frac{1}{\rho a_{\text{up}}} \text{ReLU}\left(-\frac{\partial M}{\partial z}\right). \quad (3b)$$

Here, a_{up} is the updraft area fraction, ρ is the air density, and $M = \rho a_{\text{up}} \bar{w}_{\text{up}}$ is the updraft mass flux, where \bar{w}_{up} is the updraft vertical velocity. ReLU is the rectified linear function, which ensures detrainment only occurs when the mass flux divergence is negative.

2.2.2. Nondimensionalization of Input Variables

A consequential step in designing ML problems is the choice of input variables and their preprocessing, including normalization, transformation, and feature engineering. Effective training of data-driven closures requires inputs of similar magnitude so that disproportionate importance is not assigned to variables with larger magnitudes. The online training approach complicates variable normalization since the input variables and their associated distributions are strongly dependent on entrainment mixing, and thus will vary as parameters change through the calibration process. A natural and physically motivated approach to transform input variables is to form nondimensional groups by combining dimensional variables in a manner that removes physical units. An additional advantage of doing this is that it increases the likelihood of obtaining climate-invariant closures that generalize well out of distribution (Beucler et al., 2024), in much the same way that Monin-Obukhov similarity theory is fairly generally applicable (Schneider et al., 2024).

In principle, nondimensional functions may depend on any nondimensional groups associated with lateral mixing processes. Here, nondimensional groups are found on the basis of Buckingham's Pi Theorem, which states: given N variables containing M primary dimensions, the nondimensionalized equations relating all the variables will have $(N - M)$ dimensionless groups (Buckingham, 1914). We consider a set \mathbf{D} of $N = 7$ primary variables, containing some already nondimensional quantities, namely, relative humidity (RH) and updraft area fraction (a_{up}), in addition to other variables deemed relevant for SGS turbulence and convection:

$$\mathbf{D} = \left\{ \Delta \bar{b}, \Delta \bar{w}, \overline{\text{TKE}}_{\text{env}}, z, H_{\text{scale}}, \Delta \overline{\text{RH}}, \sqrt{a_{\text{up}}} \right\}. \quad (4)$$

The set contains two length scales: the height coordinate z and the standard atmospheric scale height $H_{\text{scale}} = R_d T_{\text{ref}} / g$; $\overline{\text{TKE}}_{\text{env}}$ denotes environmental turbulent kinetic energy. Note that we use $\sqrt{a_{\text{up}}}$ instead of a_{up} because it represents a nondimensionalized length scale. Because entrainment mixing transports properties between subdomains, we defined dimensional variables as differences between the updraft and environmental properties. Using subdomain differences also ensures Galilean invariance, such that the diagnosed entrainment rates are independent of the reference frame. Given that these variables contain $M = 2$ primary dimensions (length and time), this leaves $N - M = 5$ dimensionless groups.

We use the nondimensional $\mathbf{\Pi}$ groups

$$\mathbf{\Pi} = \left\{ \frac{z \Delta \bar{b}}{\Delta \bar{w}^2}, \frac{\overline{\text{TKE}}_{\text{env}}}{\Delta \bar{w}^2}, \sqrt{a_{\text{up}}}, \Delta \overline{\text{RH}}, \frac{gz}{R_d T_{\text{ref}}} \right\}. \quad (5)$$

and refer to group i as Π_i . These $\mathbf{\Pi}$ groups, defined locally at each level of the atmosphere, serve as inputs to data-driven models that return continuous, non-negative outputs. Π_1 and Π_2 are unbounded and typically have magnitudes larger than 1, so they are normalized by characteristic values of 10^2 for Π_1 and 2 for Π_2 , such that they typically lie in the range $[-1, 1]$. Π_1 resembles the classic $\Delta \bar{b} / \Delta \bar{w}^2$ scaling introduced by Gregory (2001), and may be interpreted as a proxy for the ratio between updraft buoyancy and the updraft-environment shear. Π_2 is indicative of whether turbulent or convective kinetic energy dominate. Π_3 and Π_4 , which are already dimensionless, allow for explicitly learning the dependence of lateral mixing on updraft area and relative humidity, respectively. Finally, Π_5 serves as an easy-to-compute measure of geometric height, nondimensionalized by the density scale height.

2.2.3. Data-Driven Entrainment Architectures

The data-driven models considered for this study are linear regression and a fully connected neural network. The linear closure is a linear mapping between Π groups and the nondimensional mixing rate. A separate regression model is used for entrainment and detrainment, totaling 12 trainable mixing parameters, including bias terms. Linear regression outputs are passed through a ReLU function to ensure positivity of mixing rates. The fully connected NN contains 237 parameters with three hidden layers containing 10, 10, and 5 neurons, respectively. Neurons in all layers have ReLU activation functions. We confine ourselves here to relatively shallow network architectures, as they already yielded substantial gains in accuracy of the EDMF scheme; exploration of whether deeper networks can yield additional gains is left for future work.

2.3. GCM-Driven Simulations

We aim to learn compact representations of directly simulated, SGS processes as a function of large-scale forcings. Forcings are taken from Cloud Feedback Model Intercomparison Project sites (cfSites), which correspond to locations where high-frequency GCM output is saved for systematically diagnosing cloud feedbacks (Webb et al., 2017). To generate spread in forcings, one model from CMIP6 (CNRM-CM6) and two models from CMIP5 (HadGEM2-A and CNRM-CM5) are used, the latter two representing the upper and lower end of tropical low-cloud reflection response. The LES and EDMF scheme are driven with the same large-scale forcings from the corresponding GCM dynamical core. LES simulations are forced with GCM-prescribed tendencies for large-scale subsidence, horizontal advection, and vertical eddy advection. Additionally, entropy and total water specific humidity profiles are relaxed to the initial background GCM state with a 24 hr relaxation timescale above 3.5 km, where convective and turbulent activity cease. Momentum profiles are relaxed on a 6 hr timescale throughout the column to prevent drift. Radiation is computed interactively with RRTMG. The EDMF scheme is forced in the same manner, with the exception that radiative cooling tendencies obtained from RRTMG are prescribed from LES. LES simulations are run for 6 days; a steady state response to large-scale forcings is often observed after a couple of simulation days. Single column model simulations are ran for 3 days and more readily reach steady state. For calibration, we consider a total of 176 LES simulations across the east Pacific stratocumulus-to-cumulus transition regions. The setup discussed here is described in Z. Shen et al. (2022).

2.4. Ensemble Kalman Inversion

For calibration we employ ensemble Kalman inversion (EKI), an iterative data assimilation technique that blends Bayesian inference with stochastic ensemble sampling to efficiently find optimal parameters (M. A. Iglesias et al., 2013; Schillings & Stuart, 2017). Starting with a prior distribution over parameters, the method iteratively updates and narrows the parameter distribution by minimizing the EDMF–LES mismatch without explicitly computing gradients. After a sufficient number of iterations, the spread of the ensemble tightens around the ensemble mean, a phenomenon referred to as ensemble collapse. The method is built into a framework that optimizes EDMF parameters on the basis of LES simulations forced in the same manner. The EDMF calibration framework described here was first introduced in Lopez-Gomez et al. (2022), where further details can be found.

The Kalman update equation estimates parameters iteratively following

$$\Theta_{n+1} = \Theta_n + \text{Cov}(\Theta_n, \mathcal{G}_n) [\text{Cov}(\mathcal{G}_n, \mathcal{G}_n) + \Delta t^{-1} \Gamma]^{-1} (\mathbf{y} - \mathcal{G}_n), \quad (6)$$

where Θ is a vector containing EDMF parameters, \mathcal{G} are EDMF statistics evaluated with parameters Θ , \mathbf{y} is a vector of the reference LES statistics, and Γ is a noise covariance matrix. Subscripts denote iteration number. The sample covariance matrices $\text{Cov}(\Theta_n, \mathcal{G}_n)$ and $\text{Cov}(\mathcal{G}_n, \mathcal{G}_n)$ are computed from the ensemble members, reflecting the covariance between parameters and EDMF statistics, and within the EDMF statistics themselves. The artificial timestep is denoted Δt , and represents an EKI hyperparameter analogous to the learning rate in the gradient descent algorithm. The quantities Γ , \mathbf{y} , \mathcal{G} , and $\text{Cov}(\mathcal{G}_n, \mathcal{G}_n)$ are formed by concatenating operations over all cases in a given iteration. Statistics in \mathcal{G} and \mathbf{y} are computed with the following sequence of operations for each LES configuration. First, state variables are individually normalized by their respective time-variance over the simulation period. A time-mean is then computed over the final 12 simulation hours before a low-dimensional encoding that preserves 99% of the variance is applied through principal component projection. The projection

reduces the dimensionality of each case from 401 to 8–40. Finally, the resulting statistics are concatenated over cases to form \mathcal{G} and \mathbf{y} . The six variables whose statistics appear in the loss function are:

1. \bar{s} : entropy
2. \bar{q}_t : total water specific humidity
3. $\overline{w's'}$: vertical entropy flux
4. $\overline{w'q_t'}$: vertical total water specific humidity flux
5. \bar{q}_l : liquid water specific humidity
6. LWP: Liquid Water Path

The overbar denotes a temporal and horizontal average and primes deviations therefrom. The first five variables are vertical profiles, whereas liquid water path is a vertically integrated quantity. The pooled LES time variance, used to estimate observation noise Γ , is scaled by 0.1 for the vertical flux and liquid water specific humidity variables. We found that noise estimated from LES time variances over the full simulation results in uncertainty bands that overwhelm important details about the vertical structure of these variables. Stated differently, the temporal variability in LES simulations, used as a proxy for observation noise, likely overestimates the noise relevant for calibration for these variables. The artificial timestep Δt is determined adaptively by a Data Misfit Controller (DMC) learning rate scheduler, and generally increases with iteration number (M. Iglesias & Yang, 2021). The DMC scheduler has no hyperparameters, as timestep is computed as a function of observation noise, data misfit, and integrated timestep. The calibrations are terminated after a specified number of iterations, which are quantified below.

In the Kalman update equation, parameters encoding functional relationships of lateral mixing are denoted Θ_{ml} (machine learning parameters), and are calibrated alongside parameters Θ_p appearing in eddy diffusivity and perturbation pressure closures with imposed functional forms, which we denote physical parameters.

$$\Theta = \{\Theta_p, \Theta_{ml}\}. \quad (7)$$

Many parameter combinations lead to unstable simulations, an issue addressed by sampling from regions of the parameter space with successfully completed simulations. For a given iteration, only the subset of ensemble members with stable simulations are used to approximate the parameter distribution for the subsequent iteration, an approach detailed more in Section 3.1.1 of Lopez-Gomez et al. (2022). Model failure rates are typically 50%–80% in the initial few iterations and diminish to zero after ~ 10 iterations. To further promote stability and determine robust initial priors, we employ a 2-stage calibration process where the initial phase contains only a subset of the full LES library. The first calibration, which we denote precalibration, is performed on 5 cases using the linear regression closure and 300 ensemble members for 20 iterations. The 5 precalibration cases are representative, and span cloud regimes along the stratocumulus-to-cumulus transition. Priors for the precalibration stage are chosen from Lopez-Gomez et al. (2022) for physical parameters. Linear regression prior means are randomly drawn from a uniform distribution on the interval [0.75, 1.25] with a prior uncertainty of 5. Following this step, the neural network model is independently optimized via gradient descent to reproduce the linear regression mapping learned from EKI in the precalibration stage. For the linear closure, the second phase is initialized directly with prior means from the precalibration phase. The NN calibration is initialized with parameter means learned from gradient descent. The second phase contains all 176 LES cases and a batch size of 16 cases per iteration. Rather than evaluating the full LES data set in each iteration, 16 cases are drawn from the full data set without replacement until the entire data set is processed. A complete pass through the data set is referred to as an epoch. The final calibrations are run for 50 iterations, or ~ 4.5 epochs. The need for batching is two-fold: computational efficiency and generation of noise in the training loss. Using the full data set of 176 cases in each iteration is expensive given the runtime and memory requirements of single model runs. Additionally, variability in the forcing and cloud regimes between batches translates to variability in the evaluated loss and root mean square errors. The noise generated by the batching process inhibits convergence to local minima and is commonly used in data assimilation and machine learning (Houtekamer & Mitchell, 2001).

3. Calibration Results

3.1. Calibration Characteristics and Performance Comparison

To characterize the EKI training process, we consider the evolution of root mean squared error (rmse) separately for each of the six variables in the loss function, tracked through the final calibration and following the pre-calibration step. Figure 2 displays the evolution of rmse for the AMIP training set (left column) and a fixed set of 5 LES cases from the AMIP4K climate (right column). The AMIP4K validation cases are a representative set spanning the stratocumulus-to-cumulus transition using HadGEM2-A as the forcing model. Shading indicates the maximum and minimum rmse over ensemble members for a given iteration, as each member is associated with a unique set of parameters. A summary of rmse comparisons between the EDMF variants can be found in Appendix B. We note that the training rmse curves are noisier than the validation curves due to the batching processes. During training, the rmse for a given iteration is calculated for the 16 sampled LES cases that vary in location, season, and regime iteration-to-iteration. The validation set is intended to track generalization performance through the calibration process.

The rmse evolution represents an improvement over the precalibration posterior (full calibration prior), constrained initially by the 5 precalibration cases in the AMIP climate. Variables with larger rmse differences between the initial and final iterations benefit more from additional cases from the full AMIP training set, and vice versa. The largest differences are for \bar{q}_l and LWP, where error decreases by an order of magnitude, consistent with the sensitive and multi-scale dynamics needed to simulate cloud variables with fidelity. We note that LWP is the density weighted integral of \bar{q}_l , so the rmse values are correlated. Remaining variables, including state variables (\bar{s} , \bar{q}_l) and flux variables ($\overline{w's'}$, $\overline{w'q'_l}$), demonstrate rmse improvements of roughly 50%–75% with respect to the prior. The differences in rmse improvement may stem from observation noise differences, but these are scaled to have roughly comparable relative magnitudes, such that they hold similar weight with respect to each other in the loss. This analysis reveals that the accuracy in simulating cloud properties, through parameters that constrain \bar{q}_l , is greatly improved by expanding the number of training cases from 5 to 176.

Significant improvements of the hybrid EDMF over EDMF-20 are observed, particularly for cloud-related variables and $\overline{w's'}$. Coplotted are variable-by-variable rmse baselines evaluated with EDMF-20 over the entire AMIP data set for the training plots and the 5 AMIP4K cases in the validation plots. The most significant improvements of the hybrid EDMF over EDMF-20 are observed for \bar{q}_l , LWP, and $\overline{w's'}$. The sizable reduction of entropy flux error likely stems from the modified boundary conditions and larger entrainment rates learned near the surface. Earlier assessments of EDMF-20 demonstrated integrated entropy fluxes that were systematically biased too large, even after calibration (Lopez-Gomez, 2023). Overly warm and buoyant updrafts in EDMF-20 are likely contributors to the systematically large entropy fluxes. The updraft warm bias has been largely mitigated in the hybrid EDMF, coincident with enhanced surface entrainment that mixes cooler environmental air into the updraft and larger TKE at the surface. Less consequential improvements are identified for state variables \bar{q}_l and \bar{s} . In the validation curves, greater differences are observed between the hybrid EDMF schemes and EDMF-20, owing to data-driven closures, structural model improvements, and the larger training data set.

The comparable performance of EDMF-NN and EDMF-Linreg in training and validation metrics has several potential explanations. Differences in the learned entrainment functions are detailed further in Section 3.3. While the NN is pretrained on the linear regression model, significant prior uncertainty is introduced in the NN weights to ensure large regions of parameter space are explored beyond the linear, low-dimensional manifold. Further, given the physical structure surrounding the data-driven mixing closures, including the dimensional scale multipliers and derivation of Π groups for input, expressive and non-linear ML architectures do not appear necessary for learning the optimal mapping. The success of simple nondimensional functions may also be a consequence of simplifications made in the setup. A limitation of the training data is the use of steady large-scale forcings and LES-prescribed radiation tendencies. These preclude the simulation of high-frequency climate variability, such as the diurnal cycle of precipitation and clouds, which is more sensitive to details of entrainment (Del Genio & Wu, 2010). Nonsteady forcings with interactive radiation and deep convection cases may be needed to gain predictive benefits from more expressive mixing closures. A final contributing factor, discussed in Section 3.4, is the presence of remaining structural errors in the EDMF formulation itself, which may not be rectified through modifying the cloud mixing process.

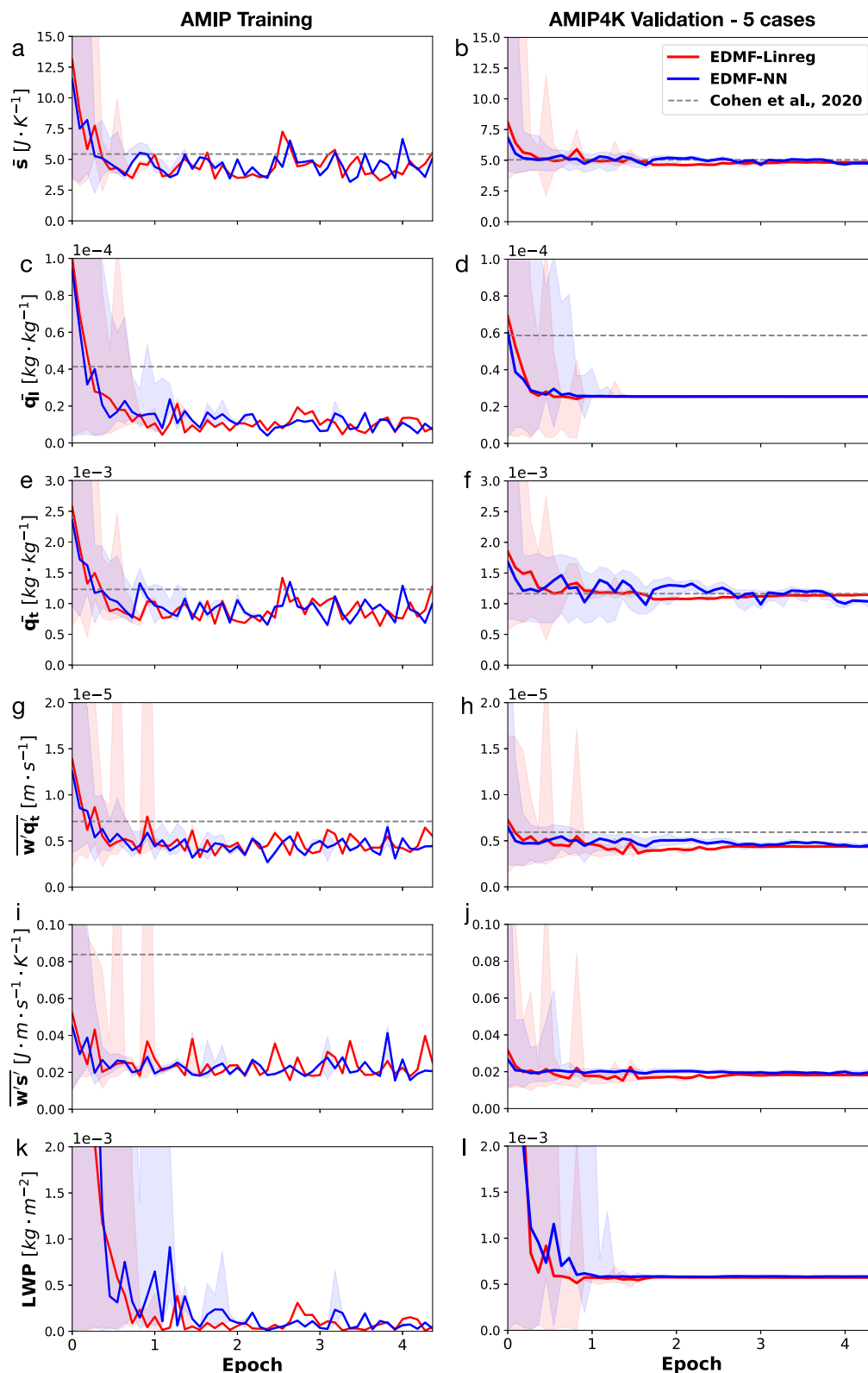


Figure 2. Root mean squared error (rmse) by variable for (left) training set from AMIP experiment and (right) validation set with five cases from the AMIP4K experiment. Shaded regions indicate min/max rmse across ensemble members for a given iteration, demonstrating ensemble spread. Dashed horizontal lines indicate baseline simulations from the EDMF-20 version described in Cohen et al. (2020). A summary of rmse comparisons can be found in Appendix B.

3.2. Generalization Performance in AMIP4K Climates

The full library of LES simulations is divided into a training and validation set on the basis of the forcing climate; the hybrid EDMF is calibrated on 176 present-day AMIP simulations and performance is evaluated on simulations from a warmer AMIP4K climate. The AMIP4K climate contains out-of-distribution large-scale forcings and surface heat fluxes. Five AMIP4K cases are chosen to track extrapolation performance through the calibration process, illustrated in the right column of Figure 2. For the chosen AMIP4K validation set, consequential performance improvements diminish after ~ 1 epoch, consistent with the training rmse. Validation rmse is noted to roughly track training rmse, with rmse for cloud-related variables \bar{q}_l and LWP containing larger extrapolation errors of $2.54 \times 10^{-5} \text{ kg} \cdot \text{kg}^{-1}$ and $5.84 \times 10^{-4} \text{ kg} \cdot \text{m}^{-2}$ for EDMF-Linreg, respectively. Nevertheless, it is found that the validation set does not enter the overfitting regime, which is characterized by a u-shaped validation curve.

Robust extrapolation performance is noted in data space as well, where key features learned in training are persistent in a simulated warmer climate. Figure 3 depicts a sampling of profiles from the AMIP4K climate across climate models, seasons, location, and cloud regimes. Optimal parameters are chosen from the ensemble member nearest to the ensemble mean at the end of the final training epoch, as the mean itself is not directly evaluated. For a given cfSite, the AMIP4K LES simulations feature changes in boundary layer depth, cloud water content, cloud depth, and vertical fluxes in response to larger surface heat fluxes and changes in local forcings due to large-scale circulation responses. Given these changes, we find hybrid EDMF simulations, trained in a cooler climate, capture these characteristics well. EDMF-20 is noted to have a large bias in \bar{q}_l near the cloud top, particularly for cumulus and transition cases. Remaining biases observed in these profiles are detailed in Section 3.4.

3.3. Learned Entrainment and Detrainment Profiles

This section turns to the assessment of learned entrainment profiles following the calibration procedure outlined above. To reiterate, the precalibration data-driven cloud mixing priors are initialized with random numbers, and closure learning is indirectly guided by the time-mean profiles alone. Focus is placed on cumulus cases, where cloud mixing is most relevant for determining the formation and behavior of clouds reliant on updraft dynamics. Figure 4 illustrates time-mean vertical profiles of the Π groups (left), nondimensional entrainment rates (middle), and total entrainment rates (right). Nonzero liquid water specific humidity (\bar{q}_l) is shaded in gray to highlight the cloud layer. The optimal parameters are chosen from the ensemble member nearest to the ensemble mean at the end of the final training epoch, as in Figure 3. The first observation to emphasize is the realism of calibrated simulations on the basis of nondimensional input groups (Figures 4a and 4d). Both EDMF-Linreg and EDMF-NN exhibit canonical characteristics of shallow convection. Notably, updraft area (Π_3) begins to shrink considerably above the cloud base due to net detrainment of mass into the environment. Near the cloud top, the updraft-environment relative humidity difference (Π_4) intensifies, where buoyant and saturated updrafts begin to penetrate into the dry, stable inversion layer. Additionally, the sub-cloud boundary layer is dominated by mixing from turbulent eddies, while the cloud layer is dominated by updraft dynamics, as indicated by the ratio of TKE to vertical velocity squared (Π_2).

The learned cloud mixing profiles themselves further demonstrate realistic and physically robust characteristics, consistent with theory surrounding lateral cloud mixing for shallow convection. Several well-established qualities of entrainment and detrainment in shallow convection include (de Rooy et al., 2013):

- A local maximum of entrainment where updrafts form;
- Net detrainment ($E - D < 0$) through much of the cloud layer;
- Strong detrainment near the cloud top, in the vicinity of a capping inversion layer.

These are consistent with theoretical work and diagnostics of lateral mixing in LES (Savre, 2022).

These key characteristics are observed in lateral mixing profiles (Figures 4c and 4f) for both EDMF-Linreg and EDMF-NN. Many SGS parameterizations feature distinct turbulent surface layer and mass-flux schemes, with the latter typically prescribing a boundary condition closure for the cloud base mass flux. Consequently, this configuration precludes both entrainment below the cloud base and strong entrainment at the cloud base. Because the EDMF scheme employed for this study is unified, updrafts may be either saturated or dry, and extended from the surface where they are generated by strong net entrainment. Coincident with near-surface updraft formation,

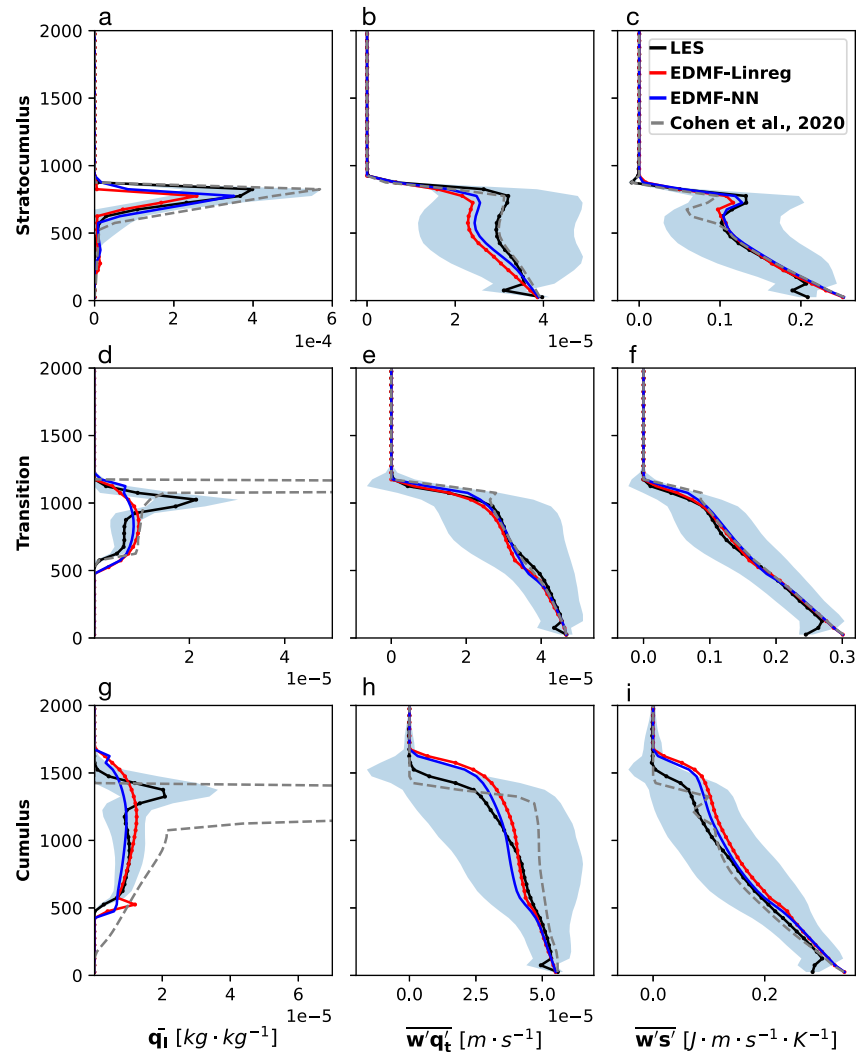


Figure 3. AMIP4K, time-mean vertical profiles of liquid water specific humidity (\bar{q}_l , left), total water specific humidity flux ($\overline{w'q_t}$, middle), and entropy flux ($\overline{w's'}$, right) from hybrid EDMF models across a sampling of climate models, seasons, geographic locations, and cloud regimes. Top row: stratoCumulus case (cfSite17) in July forced with CNRM-CM5; middle row: transition case (cfSite6) in April forced with CNRM-CM6; bottom row: cumulus case (cfSite22) in July forced with HadGEM2-A. Baseline simulations from Cohen et al. (2020) are plotted in gray dashed lines. Large-eddy simulation (LES) time-mean profiles from Z. Shen et al. (2022) are plotted in black. Calibrated EDMF simulations using a linear regression-based mixing closure (EDMF-Linreg) are depicted in red, while those with a NN-based mixing closure (EDMF-NN) are shown in blue. Light blue shading indicates the 2σ time variance, by level, from LES simulations.

large entrainment rates are observed in Figures 4c and 4f. Both closures accurately predict net detrainment above the cloud base, where entrainment rates tend to small values and detrainment grows. Finally, a global maximum in detrainment rate is observed near the cloud top.

Several core similarities and differences are discussed for the linear and NN-based entrainment closures on the basis of nondimensional rates, or the components targeted with data-driven closures. The nondimensional functions may be viewed as a multiplicative modulations of dimensional rates introduced in Equations 3a and 3b. Deviations far from unity suggest that the dimensional mixing rate does not accurately capture dynamics consistent with LES time-mean profiles. In contrast, nondimensional rates close to unity indicate that the dimensional component effectively approximates cloud mixing without need for modification. Turning to the nondimensional rates (Figures 4b and 4e), we note more consequential differences between the hybrid EDMF schemes in the detrainment rates. Notably, EDMF-NN features a secondary maximum of detrainment near the cloud base, around ~ 500 m above the surface. Such secondary local detrainment maxima are often observed in

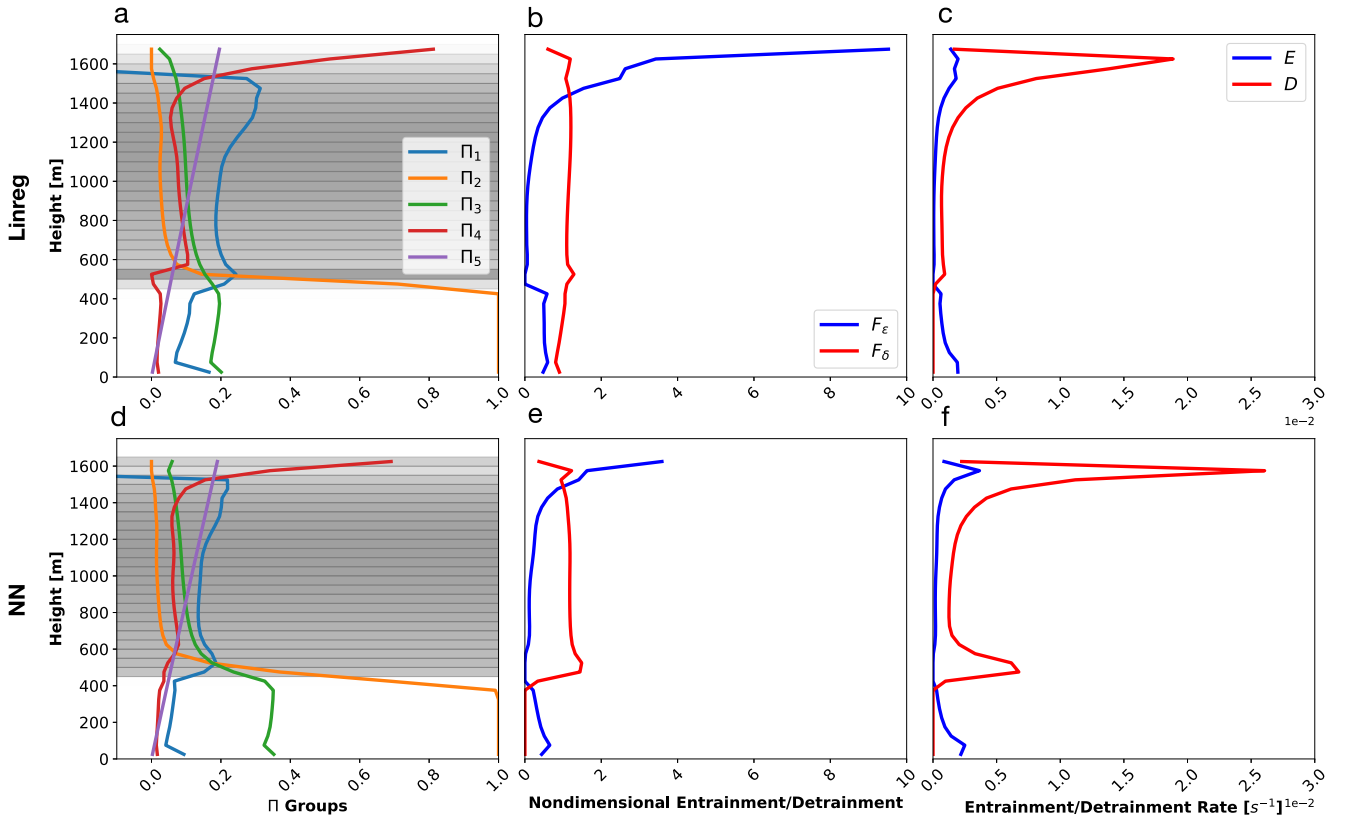


Figure 4. Time-mean vertical profiles of lateral mixing variables for cfSite22 with AMIP4K forcings, depicting shallow convection near Hawaii in July. (a, d): Nondimensional Π groups, with liquid water specific humidity (\hat{q}_l) shaded in gray. (b, e): nondimensional entrainment and detrainment (data-driven model output). (c, f): Total entrainment and detrainment rates.

LES-diagnosed detrainment rates (Romps, 2010). Generally larger detrainment rates are also observed for EDMF-NN through the cloud layer. Alternatively, EDMF-Linreg maintains a less variable nondimensional rate with height, with slight enhancement in the updraft. Focusing on nondimensional entrainment, we find stronger modulation of the dimensional scale than for detrainment. In particular, both closures demonstrate increasing modulation of the dimensional scale with height in the upper cloud levels. This indicates the $\Delta\bar{w}/z$ dimensional scale significantly underpredicts entrainment rates near the updraft top. The behavior driving this learned enhancement may surround the physical mechanisms governing cessation of updrafts, where updraft area fraction or mass flux tend to zero. Updrafts vanish by a combination of strong detrainment, which serves as a sink for area fraction, and entrainment, which diminishes upward mass flux by both reducing updraft buoyancy and entraining environmental parcels with negligible vertical momentum. Despite the two competing effects, studies point to strong net detrainment at the cloud top, as alluded to previously, which is consistent with our simulations. In the sub-cloud layer, the dimensional scale overpredicts entrainment, as indicated by nondimensional values less than unity in both schemes.

The closed-form linear expression for entrainment following the full calibration is

$$E = \frac{\Delta\bar{w}}{z} \times 6 \left[-0.05 + 0.8 \left(\frac{z\Delta\bar{b}}{\Delta\bar{w}^2} \right) + 0.6 \left(\frac{\overline{\text{TKE}}_{\text{env}}}{\Delta\bar{w}^2} \right) + -3\sqrt{a_{\text{up}}} + 3(\Delta\overline{\text{RH}}) + 0.2 \left(\frac{gz}{R_d T_{\text{ref}}} \right) \right], \quad (8)$$

and that for detrainment is

$$D = \frac{1}{\rho a_u} \text{ReLU} \left(-\frac{\partial M}{\partial z} \right) \times 8 \left[0.04 - 0.07 \left(\frac{z\Delta\bar{b}}{\Delta\bar{w}^2} \right) - 0.07 \left(\frac{\overline{\text{TKE}}_{\text{env}}}{\Delta\bar{w}^2} \right) + 0.8\sqrt{a_{\text{up}}} - 0.2(\Delta\overline{\text{RH}}) + 0.5 \left(\frac{gz}{R_d T_{\text{ref}}} \right) \right]. \quad (9)$$

These are determined from the ensemble member nearest to the mean in the final training epoch. These functional relationships may be used to understand the vertical structure of nondimensional mixing in the context of Figure 4. In the sub-cloud surface layer, where a local entrainment maximum is observed (Figures 4c and 4f), the linear model has strong contributions from Π_2 as a consequence of large TKE. Above the surface layer, the increase of nondimensional entrainment with height has large contributions from gradually decreasing area fraction (Π_3) through the cloud layer and sharply increasing updraft-environment relative humidity difference (Π_4) near the cloud top (Figures 4a and 4d). The linear nondimensional detrainment rates demonstrate weaker variation with height. The Π groups themselves contain covariances, so variable importance cannot not be read off explicitly from Equation 8 and 9. Because the full calibration is initialized with parameter means from precalibration, differences in the final parameter values indicate sensitivity to number of training cases, particularly when going from 5 to 176 cases. We find the training data sensitivity to be parameter-dependent. The entrainment weights for Π_3 and Π_4 , in particular, demonstrate the most sensitivity. For entrainment, the full calibration modified the Π_3 weight by a factor of ~ 2 and the Π_4 weight by a factor of ~ 3 . Alternatively, the detrainment parameters for Π_3 and the bias have little sensitivity beyond 5 cases, and are modified by $<10\%$ in the full calibration. Remaining parameters exhibit intermediate sensitivities. In Appendix C, we provide a comparison of the final linear regression weights following each experiments, as well as precalibration uncertainty estimates from the Calibrate, Emulate, Sample framework (Cleary et al., 2021).

3.4. Beyond Calibration: Addressing Structural Errors

Post-calibration, persisting discrepancies between the LES and EDMF may be attributed to three primary contributions: the EKI optimizer, the inverse problem setup, and inherent biases in the underlying physical forward model or data, in this case, the structure and assumptions of the EDMF scheme. The performance of the EKI optimizer, as determined by its convergence, may be sensitive to EKI settings and hyperparameters. Among the most consequential choices are the EKI artificial timestepper and the batch size. Sensitivity to constant artificial timestep values in previous work (Lopez-Gomez et al., 2022) is addressed here by using a hyperparameter-free adaptive timestep (DMC) that increases through the calibration process. For batching, we chose the largest batch size feasible given computational limitations. It is found that batch sizes smaller than ~ 10 generate excessive noise in the loss, preventing descent of the ensemble mean to lower values and convergence of the EKI algorithm. Additional biases may persist as a result of the problem setup, such as the input variables selected for data-driven closures and the choice of priors. In addition to addressing instabilities, the precalibration procedure reduces sensitivities to the priors. Precalibration is initialized with large prior uncertainties over parameters with a relatively large number of ensemble members (300), allowing broad exploration of the parameter space and narrowing of the posterior on the basis of a small but representative data set. While these approaches curtail EDMF-LES discrepancies and mitigate convergence to local minima, it is possible that more advanced strategies are needed to initialize, pretrain, and calibrate the NN-based EDMF. Attempts to initiate the EDMF-NN calibrations directly with Xavier initialization (Glorot & Bengio, 2010) produced EKI calibrations that exhibited high ensemble failure rates and minimal convergence of the loss function across a range of prior uncertainties.

Structural error denotes errors arising from the design of the EDMF scheme itself, including but not limited to the formulation of other closures, boundary conditions, and assumptions made in deriving the EDMF equations. Such limitations may not be corrected by calibration, but must be addressed by modifying the anatomy of the EDMF scheme or adding structural error models within the governing EDMF equations. Relative to Lopez-Gomez et al. (2022), this study addressed three structural errors by modifying the EDMF equations and boundary conditions:

1. A strong warm bias near the surface, resulting from a TKE minimum in the bottom cell center, addressed by implementing a bottom flux boundary condition for the TKE equation;
2. Calibrations with near-zero entrainment throughout the vertical profile, addressed by implementing a free boundary condition on updraft area in the bottom cell center;
3. Divergence of area fraction to values close to one, addressed by choosing a dimensional scale for detrainment that ensures area fraction gradually tends to zero when the mass flux gradient is negative.

These modifications led to both improved training and validation errors as well as more realistic cloud mixing profiles following calibration.

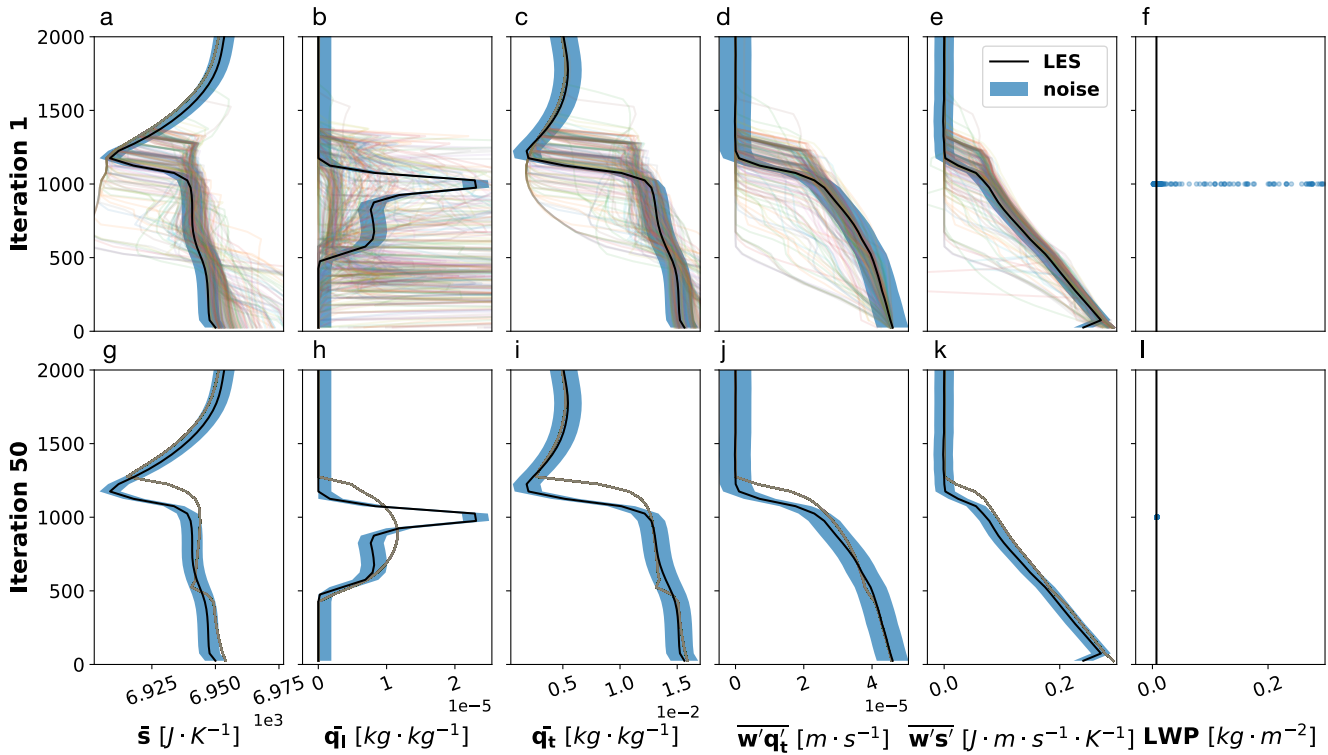


Figure 5. Ensemble spread of EDMF-Linreg for all loss function variables in (top) first iteration and (bottom) final iteration. Large-eddy simulation (LES) time-mean profiles are plotted in black (Z. Shen et al., 2022), and each colored lines represents the evaluation from an ensemble member. Blue shading indicates the 2σ observation noise used by EKI, calculated from the pooled variance across levels in LES simulations.

Remaining structural errors primarily involve biases in the depth of the mixed layer and cloud top \bar{q}_i maxima. First, we note an underestimation of capping stratocumulus clouds in stratocumulus-topped cumulus forcing regimes, as demonstrated by \bar{q}_i profiles in Figures 3d and 5h. While relatively low \bar{q}_i errors are observed for layers composed of cumulus clouds in these regimes, below roughly 1,000 m in Figure 3d and 800 m in Figure 5h, the grid-mean \bar{q}_i is biased systematically low at cloud tops. Transition cases demonstrating this bias contain saturated updrafts in the cloud layer, but fail to saturate the environment at the level stratocumulus clouds are observed in LES simulations. Because stratocumulus dynamics are dominated by environmental mixing, rather than updraft dynamics, this likely indicates a bias in the TKE equations or other environmental factors. This hypothesis is further supported by the initial spread of \bar{q}_i profiles across ensemble members in data space, illustrated in Figure 5b. The initial iteration contains sizable spread in parameter values, consistent with the prior, and is indicative of the data space subsequent iterations will explore. Characteristics, such as capping stratocumulus clouds, not loosely demonstrated by ensemble members during the initial iterations are unlikely to be developed in later iterations, implying a systematic bias in the model or prior means that are far removed from the optimal solution for a given case. We found the bias to be persistent across many calibration in offline experiments varying the precalibration set and EKI settings. The bias is further demonstrated by systematic collapse of ensembles in the final iteration far beyond the envelope of observation noise (Figure 5h). Cloud top maxima of \bar{q}_i are also observed for LES simulations of pure shallow convection, but these features may be an artifact of microphysics in LES simulations. Anvil-like structures in the LES shallow convection cases are coincident with vertical maxima of cloud fraction, and may not be desirable to fit to.

Second, we note a bias in mixed layer depth for some cases, resulting in biases across variables near the cloud top. This is evident in the shallow cumulus case illustrated in Figure 3, where the mixed layer becomes ~ 100 m too deep, as evidenced by the vertical fluxes in panels h, i. As a consequence, the cloud also develops too deeply (Figure 3g). While most cases capture the depth of the mixed layer with high fidelity, cases with the most prominent bias in cloud top stratocumulus structures tend to coincide with a bias in the mixed layer depth. Remaining structural errors may be rectified in future work by replacing additional closures with data-driven

models or learning structural error models as additional additive terms that modify EDMF tendency equations (Wu et al., 2023). With the latter strategy, care must be taken to ensure conservation of mass, momentum, and energy. Given biases in the depth of the mixed layer and cloud top stratocumulus structures in transition cases, we believe adding data-driven closures or error models to the TKE equation would help address these issues.

4. Concluding Remarks

In this study, our aim was to develop realistic hybrid SGS models that combine generalizability with interpretability, targeting the challenging Pacific stratocumulus-to-cumulus transition—a region notorious for being particularly error-prone in state-of-the-art climate models. The primary contribution of this paper is the demonstration of online learning of a 1D hybrid model in realistic climate settings, a step needed to eventually apply such methods in operational GCMs. Application in realistic setups may require pretraining more expressive data-driven components (NNs) to obtain sensible priors, failure handling mechanisms to address numerically unstable simulations in the training process, and procedures or guidelines for identifying remaining structural biases. Development of hybrid models benefits from a bidirectional workflow, where online learning is informative about where structural model biases might lie, and calibrations of data-driven components help improve the predictive power of hybrid models. Finally, and critical in the development of hybrid SGS models, is the assessment of physical validity alongside predictive power. Success of the hybrid EDMF is particularly evident in the realism of cloud mixing closures, which were learned indirectly from extensive LES data with no direct prior information about entrainment and detrainment. The learned closures align closely with existing theoretical understanding and LES-diagnosed characteristics of lateral cloud mixing as it relates to convective and cloud dynamics, reinforcing the model's scientific validity. Furthermore, our results highlight the hybrid model's predictive power, with substantial improvements over a baseline EDMF tuned to match field campaigns. We observe that performance improvements translate to an out-of-distribution AMIP4K climate, as assessed by rmse and qualitative analysis of physical profiles. This generalizability is crucial for the model's application to prediction of future climate scenarios in GCMs.

The online learning approach for hybrid modeling presents several advantages over offline, fully data driven alternatives. The EKI framework allows for indirectly training SGS model components on the basis of observable statistics or quantities appropriate for long-term climate model projections. While the study focused on high-resolution simulations for training, this may be extended to include sparse observations in the loss function. Numerical instabilities resulting from unstable parameter combinations are directly addressed in the training process, reducing the likelihood of instabilities when the parameterization is incorporated in operational GCMs. Additionally, data-driven components of a hybrid model can be more easily isolated and reasoned about, giving stronger confidence in out-of-distribution predictions of future climate states and promoting physical process understanding.

Despite these promising developments, there are remaining avenues for improving the hybrid EDMF scheme. The paper highlights that the reliance on steady large-scale forcings and prescribed radiation tendencies in the training data limits the ability to learn phenomena important for capturing high-frequency climate variability, such as the diurnal cycle. Additional data sets of high-resolution simulations, such as those introduced by Chammas et al. (2023) and Yu et al. (2024), would likely improve performance over a broader range of forcings and atmospheric regimes. Additionally, some errors in the structure of the model persist after calibration, resulting in a form of underfitting. Remaining structural errors may be remedied in future work by replacing additional closures with expressive, data-driven components or learning structural error corrections as additional additive terms that modify EDMF tendency equations. One avenue is to target closures in the environmental TKE equation, as the data-driven lateral mixing closures presented here primarily affect updraft characteristics and mass flux. Because our EDMF scheme uses a 1.5-order Mellor-Yamada turbulence closure, a natural target is the mixing length function, which determines environmental turbulent diffusivity and viscosity (Mellor & Yamada, 1982). Future work should focus on these aspects, in addition to more expansive training data sets, to ensure that the hybrid modeling approach can be effectively applied in operational Earth system models.

Appendix A: Hybrid EDMF Bottom Boundary Conditions

A1. Updraft Area

The inhomogeneous Dirichlet boundary condition on area in EDMF-20 is replaced by a free boundary condition, where updraft area is generated directly by entrainment and detrainment source terms at the bottom boundary. Because area is a prognostic variable in the EDMF equations, choices must be made about how the boundary conditions are specified. The EDMF continuity equation for a single updraft reads

$$\frac{\partial(\rho a_{\text{up}})}{\partial t} = -\nabla_h \cdot (\rho a_{\text{up}} \langle u_h \rangle) - \frac{\partial(\rho a_{\text{up}} \bar{w}_{\text{up}})}{\partial z} + \rho a_{\text{up}}(E - D) \quad (\text{A1})$$

where $\langle u_h \rangle$ is the average grid-scale horizontal velocity, ∇_h is the horizontal divergence, a_{up} is the updraft area fraction, \bar{w}_{up} is the updraft vertical velocity, ρ is the air density, and E and D are entrainment and detrainment, respectively.

The bottom area fraction was previously specified as an EDMF parameter a_s , typically chosen as 0.1, which remained fixed in all simulations (Cohen et al., 2020; Lopez-Gomez et al., 2022; Tan et al., 2018). The Dirichlet boundary condition on area was defined as

$$\rho a(z_0) = \rho a_s \quad (\text{A2})$$

where z_0 is the height of the interior point adjacent to the bottom boundary. Removing the surface area parameter and allowing for a free boundary condition permits the generation of surface-based updrafts directly from source terms. The modification allows updrafts to be generated by net entrainment ($E - D > 0$) or grid-scale horizontal convergence near the surface, and thus vary with environmental conditions.

A2. Turbulent Kinetic Energy

We substitute the TKE Dirichlet boundary condition in EDMF-20 by a flux boundary condition at the bottom boundary. The Dirichlet boundary condition was formulated as

$$\overline{\text{TKE}}_{\text{env}}(z_0) = \kappa_*^2 u_*^2 \quad (\text{A3})$$

where $\overline{\text{TKE}}_{\text{env}}$ represents the environmental TKE, κ_* is the ratio of rms turbulent velocity to the friction velocity (an EDMF parameter), u_* is the friction velocity, and z_0 is the height of the interior point adjacent to the boundary.

We replaced this formulation by a flux boundary condition on the TKE flux at the bottom boundary. To obtain the flux boundary condition, the following simplifying assumptions are made:

1. The mixing length in the surface layer is limited by the distance to the boundary.
2. Storage and mean advection of $\overline{\text{TKE}}_{\text{env}}$ are neglected. This is a good approximation in the surface layer, where TKE is roughly constant.
3. Horizontal derivatives are small compared to the vertical derivatives close to the boundary (the boundary layer approximation).
4. The velocity-pressure gradient correlation term can be neglected. This assumption is consistent with the impenetrability condition for the subdomains and the closure for perturbation pressure in the EDMF model.

These approximations lead to the flux-gradient relation at the surface

$$\rho a_{\text{env}} \overline{w_0' \text{TKE}_{\text{env}}'} \Big|_{z_0} = \rho a_{\text{env}} (1 - c_d c_m \kappa_*^4) u_*^2 \|u_{p,\text{int}}\|, \quad (\text{A4})$$

where a_{env} is the environmental area fraction, $u_{p,\text{int}}$ is the near-surface velocity component parallel to the surface, c_d is the turbulent dissipation coefficient, and c_m is the eddy viscosity coefficient (Lopez-Gomez et al., 2022). The modification allows the surface TKE to vary more strongly with environmental conditions.

Appendix B: RMSE Tables

Table B1

EDMF version—AMIP	\bar{s}	\bar{q}_l	\bar{q}_t	$\bar{w}^T \bar{q}_t$	$\bar{w}^T \bar{s}^T$	LWP
EDMF-NN	5.55	8.26e−06	1.29e−03	5.54e−06	2.54e−02	4.72e−05
EDMF-Linreg	5.10	7.25e−06	1.00e−03	4.45e−06	2.06e−02	3.14e−05
Cohen et al., 2020	5.43	4.13e−05	1.23e−03	7.12e−06	8.38e−02	1.79e−01

Note. Reported RMSE Values for EDMF-NN and EDMF-Linreg are the Ensemble-Averaged RMSE in the Final Iteration.

Table B2

EDMF version—AMIP4K	\bar{s}	\bar{q}_l	\bar{q}_t	$\bar{w}^T \bar{q}_t$	$\bar{w}^T \bar{s}^T$	LWP
EDMF-NN	4.84	2.54e−05	1.14e−03	4.37e−06	1.82e−02	5.73e−04
EDMF-Linreg	4.78	2.54e−05	1.06e−03	4.44e−06	1.88e−02	5.84e−04
Cohen et al., 2020	5.03	5.86e−05	1.16e−03	5.93e−06	7.93e−01	2.13e−01

Appendix C: Parameter Sensitivities

This analysis compares linear regression parameters for the 5-case precalibration (precal) and 176-case full calibration (full cal), including precalibration uncertainty estimates using the Calibrate, Emulate, Sample (CES) framework (Cleary et al., 2021). The precal prior and posterior distributions help contextualize shifts in the final weights between experiments. Figure C1 shows the precal prior and posterior distributions, coplotted with the final parameter values for each experiment. The linear weights multiplying Π_i are labeled C_i^e for entrainment and C_i^d for detrainment. The corresponding bias terms are labeled bias^e and bias^d for entrainment and detrainment, respectively. We find the training data sensitivity to be parameter-dependent, as indicated by varying degrees of modification to the final parameter values between experiments. In precalibration, the emulation step consists of training Gaussian processes containing a radial basis function kernel on parameter-to-data pairs from the ensemble. We train the emulator on the first iteration with a failure rate below 50% (iteration 4), and include iterations 8 and 16 to better emulate regions of the parameter space where the ensemble is converging. The sample step probes uncertainty via Markov Chain Monte Carlo (MCMC) sampling of the parameter space around the precal final mean parameter values. We use 100,000 samples for MCMC, with the first 2,000 samples discarded as burn-in to ensure the chain reaches equilibrium and mitigate the impact of initialization bias.

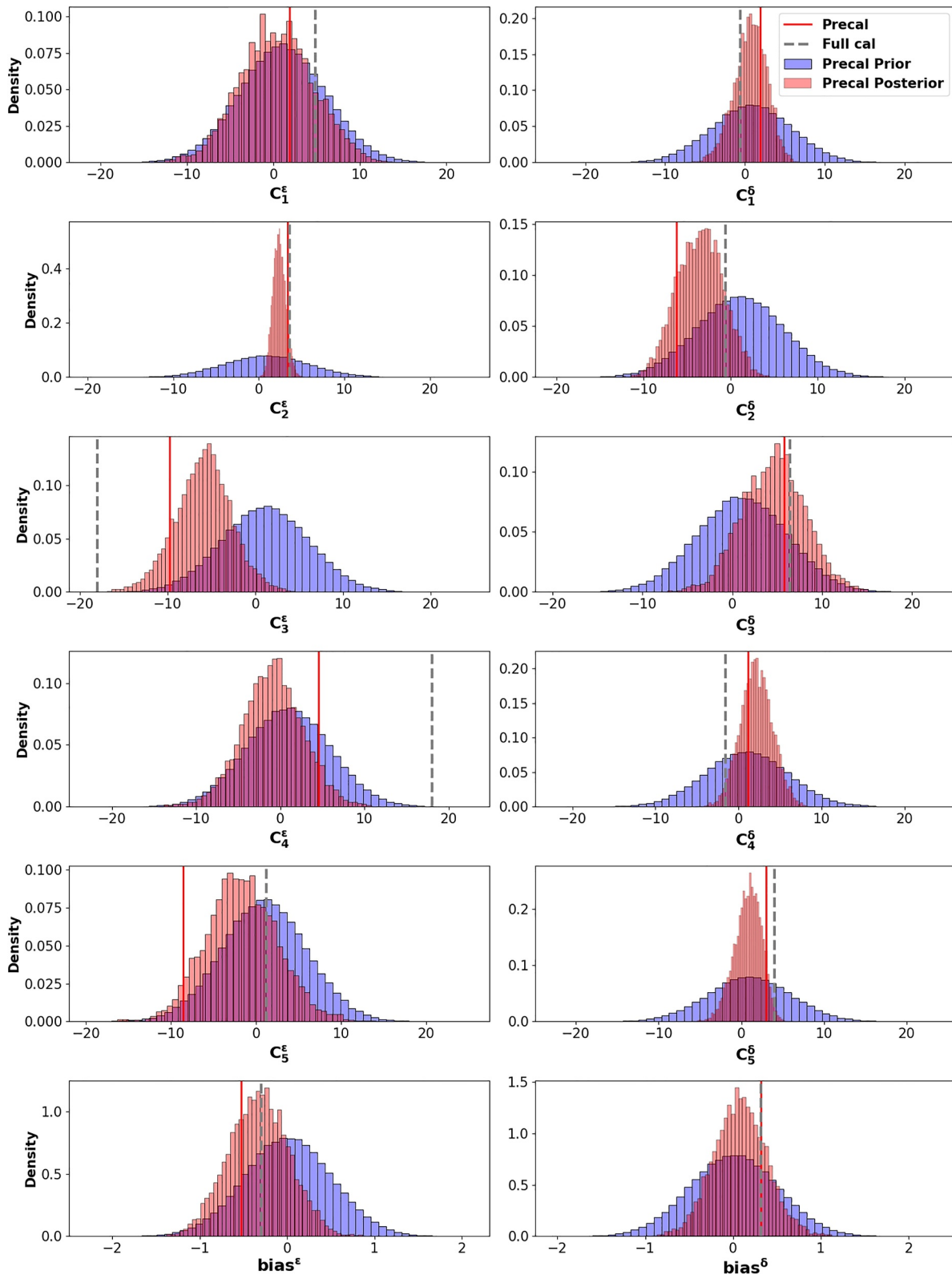


Figure C1. Prior and posterior parameter uncertainty estimated by Calibrate, Emulate, Sample (Cleary et al., 2021) for the 5-case precalibration. Blue distributions indicate the prior and red distributions indicate the posterior. Vertical lines mark final parameter values for the precalibration (precal) in solid red and the 176-case full calibration (full cal) in dashed gray, determined by taking the ensemble member nearest to the ensemble mean in the final iteration. Entrainment parameters are in the left column and detrainment parameters are in the right.

Data Availability Statement

The calibration pipeline and underlying EDMF model used for this work are available as open-source Julia packages. The EDMF single column model is TurbulenceConvection.jl v1.3.6, available at <https://doi.org/10.5281/zenodo.13733436> (Kawczynski et al., 2024). The calibration pipeline for the EDMF is implemented in CalibrateEDMF.jl v0.8.1 (<https://doi.org/10.5281/zenodo.13738494>) (Lopez-Gomez et al., 2024), and the underlying ensemble Kalman inversion algorithms are part of EnsembleKalmanProcesses.jl v1.1.5 (<https://doi.org/10.5281/zenodo.10146103>) (Dunbar et al., 2023). Visualization tools for calibration results are available alongside the calibration data at <https://doi.org/10.5281/zenodo.13743167> (Christopoulos, 2024). The PyCLES large-eddy simulation output used for calibration is available on CaltechDATA (Z. Shen, 2022).

Acknowledgments

We thank Zhaoyi Shen, Anna Jaruga, and Haakon Ervik for significant contributions to the development of the EDMF, which is the basis of this calibration work. This research was supported by Schmidt Sciences, LLC, by the U.S. National Science Foundation (Grant number AGS-1835860), and by the Office of Naval Research (Grant number N00014-23-1-2654). Tom Beucler acknowledges partial funding from the Swiss State Secretariat for Education, Research and Innovation (SERI) for the Horizon Europe project AI4PEX (Grant agreement ID: 101137682).

References

- Ackerman, A. S., VanZanten, M. C., Stevens, B., Savic-Jovicic, V., Bretherton, C. S., Chlond, A., et al. (2009). Large-eddy simulations of a drizzling, stratocumulus-topped marine boundary layer. *Monthly Weather Review*, *137*(3), 1083–1110. <https://doi.org/10.1175/2008MWR2582.1>
- Beucler, T., Gentine, P., Yuval, J., Gupta, A., Peng, L., Lin, J., et al. (2024). Climate-invariant machine learning. *Science Advances*, *10*(6), eadj7250. <https://doi.org/10.1126/sciadv.adj7250>
- Bony, S., Stevens, B., Frierson, D. M. W., Jakob, C., Kageyama, M., Pincus, R., et al. (2015). Clouds, circulation and climate sensitivity. *Nature Geoscience*, *8*(4), 261–268. <https://doi.org/10.1038/ngeo2398>
- Boutle, I. A., Eyre, J. E. J., & Lock, A. P. (2014). Seamless stratocumulus simulation across the turbulent gray zone. *Monthly Weather Review*, *142*(4), 1655–1668. <https://doi.org/10.1175/MWR-D-13-00229.1>
- Brajard, J., Carrasi, A., Bocquet, M., & Bertino, L. (2021). Combining data assimilation and machine learning to infer unresolved scale parametrization. *Philosophical Transactions of the Royal Society A: Mathematical, Physical & Engineering Sciences*, *379*(2194), 20200086. <https://doi.org/10.1098/rsta.2020.0086>
- Brenowitz, N. D., Beucler, T., Pritchard, M., & Bretherton, C. S. (2020). Interpreting and stabilizing machine-learning parametrizations of convection. *Journal of the Atmospheric Sciences*, *77*(12), 4357–4375. <https://doi.org/10.1175/JAS-D-20-0082.1>
- Brenowitz, N. D., & Bretherton, C. S. (2019). Spatially extended tests of a neural network parametrization trained by coarse-graining. *Journal of Advances in Modeling Earth Systems*, *11*(8), 2728–2744. <https://doi.org/10.1029/2019ms001711>
- Brient, F., & Schneider, T. (2016). Constraints on climate sensitivity from space-based measurements of low-cloud reflection. *Journal of Climate*, *29*(16), 5821–5835. <https://doi.org/10.1175/JCLI-D-15-0897.1>
- Brown, A. R., Cederwall, R. T., Chlond, A., Duynkerke, P. G., Golaz, J. C., Khairoutdinov, M., et al. (2002). Large-eddy simulation of the diurnal cycle of shallow cumulus convection over land. *Quarterly Journal of the Royal Meteorological Society*, *128*(582), 1075–1093. <https://doi.org/10.1256/003590002320373210>
- Buckingham, E. (1914). On physically similar systems; illustrations of the use of dimensional equations. *Physical Review*, *4*(4), 345–376. <https://doi.org/10.1103/PhysRev.4.345>
- Chammas, S., Wang, Q., Schneider, T., Ihme, M., Chen, Y., & Anderson, J. (2023). Accelerating large-eddy simulations of clouds with tensor processing units. *Journal of Advances in Modeling Earth Systems*, *15*(10), e2023MS003619. <https://doi.org/10.1029/2023MS003619>
- Christopoulos, C. (2024). Data for “online learning of entrainment closures in a hybrid machine learning parameterization” [Collection]. *Zenodo*. <https://doi.org/10.5281/zenodo.13743167>
- Cleary, E., Garbuno-Inigo, A., Lan, S., Schneider, T., & Stuart, A. M. (2021). Calibrate, emulate, sample. *Journal of Computational Physics*, *424*, 109716. <https://doi.org/10.1016/j.jcp.2020.109716>
- Cohen, Y., Lopez-Gomez, I., Jaruga, A., He, J., Kaul, C. M., & Schneider, T. (2020). Unified entrainment and detrainment closures for extended eddy-diffusivity mass-flux schemes. *Journal of Advances in Modeling Earth Systems*, *12*(9), e2020MS002162. <https://doi.org/10.1029/2020MS002162>
- Črnivec, N., Cesana, G., & Pincus, R. (2023). Evaluating the representation of tropical stratocumulus and shallow cumulus clouds as well as their radiative effects in CMIP6 models using satellite observations. *Journal of Geophysical Research: Atmospheres*, *128*(23), e2022JD038437. <https://doi.org/10.1029/2022JD038437>
- Del Genio, A. D., & Wu, J. (2010). The role of entrainment in the diurnal cycle of continental convection. *Journal of Climate*, *23*(10), 2722–2738. <https://doi.org/10.1175/2009JCLI3340.1>
- de Rooy, W. C., Bechtold, P., Fröhlich, K., Hohenegger, C., Jonker, H., Mironov, D., et al. (2013). Entrainment and detrainment in cumulus convection: An overview. *Quarterly Journal of the Royal Meteorological Society*, *139*(670), 1–19. <https://doi.org/10.1002/qj.1959>
- de Rooy, W. C., & Pier Siebesma, A. (2010). Analytical expressions for entrainment and detrainment in cumulus convection: Analytical Expressions for Entrainment and Detrainment. *Quarterly Journal of the Royal Meteorological Society*, *136*(650), 1216–1227. <https://doi.org/10.1002/qj.640>
- Dunbar, O. R. A., Constantinou, N. C., Lopez-Gomez, I., Garbuno-Inigo, A., Bolewski, J., Bach, E., et al. (2023). EnsembleKalmanProcesses.jl (version 1.1.5) [Software]. *Zenodo*. <https://doi.org/10.5281/zenodo.10146103>
- Dunbar, O. R. A., Duncan, A. B., Stuart, A. M., & Wolfram, M.-T. (2022). Ensemble inference methods for models with noisy and expensive likelihoods. *SIAM Journal on Applied Dynamical Systems*, *21*(2), 1539–1572. <https://doi.org/10.1137/21M1410853>
- Dunbar, O. R. A., Garbuno-Inigo, A., Schneider, T., & Stuart, A. M. (2021). Calibration and uncertainty quantification of convective parameters in an idealized GCM. *Journal of Advances in Modeling Earth Systems*, *13*(9), e2020MS002454. <https://doi.org/10.1029/2020MS002454>
- Frezat, H., Le Sommer, J., Fablet, R., Balarac, G., & Lguensat, R. (2022). A posteriori learning for quasi-geostrophic turbulence parametrization. *Journal of Advances in Modeling Earth Systems*, *14*(11), e2022MS003124. <https://doi.org/10.1029/2022MS003124>
- Glorot, X., & Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics* (Vol. 9, 249–256). JMLR Workshop and Conference Proceedings. <https://proceedings.mlr.press/v9/glorot10a>
- Grabowski, W. W., Bechtold, P., Cheng, A., Forbes, R., Halliwell, C., Khairoutdinov, M., et al. (2006). Daytime convective development over land: A model intercomparison based on LBA observations. *Quarterly Journal of the Royal Meteorological Society*, *132*(615), 317–344. <https://doi.org/10.1256/qj.04.147>

- Gregory, D. (2001). Estimation of entrainment rate in simple models of convective clouds. *Quarterly Journal of the Royal Meteorological Society*, 127(571), 53–72. <https://doi.org/10.1002/qj.49712757104>
- Griewank, P. J., Heus, T., & Neggers, R. A. J. (2022). Size-dependent characteristics of surface-rooted three-dimensional convective objects in continental shallow cumulus simulations. *Journal of Advances in Modeling Earth Systems*, 14(3), e2021MS002612. <https://doi.org/10.1029/2021MS002612>
- Heinze, R., Mironov, D., & Raasch, S. (2015). Second-moment budgets in cloud topped boundary layers: A large-eddy simulation study. *Journal of Advances in Modeling Earth Systems*, 7(2), 510–536. <https://doi.org/10.1002/2014MS000376>
- Holland, J. Z., & Rasmusson, E. M. (1973). Measurements of the atmospheric mass, energy, and momentum budgets over a 500-kilometer square of tropical ocean. *Monthly Weather Review*, 101(1), 44–55. [https://doi.org/10.1175/1520-0493\(1973\)101<0044:MOTAME>2.3.CO;2](https://doi.org/10.1175/1520-0493(1973)101<0044:MOTAME>2.3.CO;2)
- Houtekamer, P. L., & Mitchell, H. L. (2001). A sequential ensemble Kalman filter for atmospheric data assimilation. *Monthly Weather Review*, 129(1), 123–137. [https://doi.org/10.1175/1520-0493\(2001\)129<0123:ASEKFF>2.0.CO;2](https://doi.org/10.1175/1520-0493(2001)129<0123:ASEKFF>2.0.CO;2)
- Huang, D. Z., Huang, J., Reich, S., & Stuart, A. M. (2022a). Efficient derivative-free Bayesian inference for large-scale inverse problems. *Inverse Problems*, 38(12), 125006. (arXiv:2204.04386 [cs, math]). <https://doi.org/10.1088/1361-6420/ac99fa>
- Huang, D. Z., Schneider, T., & Stuart, A. M. (2022b). Iterated Kalman methodology for inverse problems. *Journal of Computational Physics*, 463, 111262. <https://doi.org/10.1016/j.jcp.2022.111262>
- Iglesias, M., & Yang, Y. (2021). Adaptive regularisation for ensemble Kalman inversion. *Inverse Problems*, 37(2), 025008. <https://doi.org/10.1088/1361-6420/abd29b>
- Iglesias, M. A., Law, K. J. H., & Stuart, A. M. (2013). Ensemble Kalman methods for inverse problems. *Inverse Problems*, 29(4), 045001. <https://doi.org/10.1088/0266-5611/29/4/045001>
- Kawczynski, C., Jaruga, A., Lopez-Gomez, I., Christopoulos, C., Ervik, H. L. L., Shen, Z., & Nelsen, N. (2024). TurbulenceConvection.jl (version 1.3.6) [Software]. <https://doi.org/10.5281/zenodo.13733436>
- Kochkov, D., Smith, J. A., Alieva, A., Wang, Q., Brenner, M. P., & Hoyer, S. (2021). Machine learning–accelerated computational fluid dynamics. *Proceedings of the National Academy of Sciences*, 118(21), e2101784118. <https://doi.org/10.1073/pnas.2101784118>
- Kochkov, D., Yuval, J., Langmore, I., Norgaard, P., Smith, J., Mooers, G., et al. (2024). Neural general circulation models for weather and climate. *Nature*, 632(8027), 1060–1066. <https://doi.org/10.1038/s41586-024-07744-y>
- Kovachki, N. B., & Stuart, A. M. (2019). Ensemble kalman inversion: A derivative-free technique for machine learning tasks. *Inverse Problems*, 35(9), 095005. (arXiv: 1808.03620). <https://doi.org/10.1088/1361-6420/ab1c3a>
- Krasnopolsky, V. M., Fox-Rabinovitz, M. S., & Belochitski, A. A. (2013). Using ensemble of neural networks to learn stochastic convection parameterizations for climate and numerical weather prediction models from data simulated by a cloud resolving model. *Advances in Artificial Neural Systems*, 2013, 1–13. <https://doi.org/10.1155/2013/485913>
- List, B., Chen, L.-W., & Thuerey, N. (2022). Learned turbulence modelling with differentiable fluid solvers: Physics-based loss functions and optimisation horizons. *Journal of Fluid Mechanics*, 949, A25. <https://doi.org/10.1017/jfm.2022.738>
- Lopez-Gomez, I. (2023). *A unified data-informed model of turbulence and convection for climate prediction*. Doctoral dissertation. California Institute of Technology. <https://thesis.library.caltech.edu/15063/>
- Lopez-Gomez, I., Christopoulos, C., Ervik, H. L. L., Kawczynski, C., Jaruga, A., & Shen, Z. (2024). CalibrateEDMF.jl (version 0.8.1) [Software]. <https://doi.org/10.5281/zenodo.13738494>
- Lopez-Gomez, I., Christopoulos, C., Langeland Ervik, H. L., Dunbar, O. R. A., Cohen, Y., & Schneider, T. (2022). Training physics-based machine-learning parameterizations with gradient-free ensemble Kalman methods. *Journal of Advances in Modeling Earth Systems*, 14(8), e2022MS003105. <https://doi.org/10.1029/2022MS003105>
- Lopez-Gomez, I., Cohen, Y., He, J., Jaruga, A., & Schneider, T. (2020). A generalized mixing length closure for eddy-diffusivity mass-flux schemes of turbulence and convection. *Journal of Advances in Modeling Earth Systems*, 12(11), e2020MS002161. <https://doi.org/10.1029/2020MS002161>
- MacArt, J. F., Sirignano, J., & Freund, J. B. (2021). Embedded training of neural-network subgrid-scale turbulence models. *Physical Review Fluids*, 6(5), 050502. <https://doi.org/10.1103/PhysRevFluids.6.050502>
- Meehl, G. A., Senior, C. A., Eyring, V., Flato, G., Lamarque, J.-F., Stouffer, R. J., et al. (2020). Context for interpreting equilibrium climate sensitivity and transient climate response from the CMIP6 Earth system models. *Science Advances*, 6(26), eaba1981. <https://doi.org/10.1126/sciadv.aba1981>
- Mellor, G. L., & Yamada, T. (1982). Development of a turbulence closure model for geophysical fluid problems. *Reviews of Geophysics*, 20(4), 851–875. <https://doi.org/10.1029/RG020i004p00851>
- Myers, T. A., Scott, R. C., Zelinka, M. D., Klein, S. A., Norris, J. R., & Caldwell, P. M. (2021). Observational constraints on low cloud feedback reduce uncertainty of climate sensitivity. *Nature Climate Change*, 11(6), 501–507. <https://doi.org/10.1038/s41558-021-01039-0>
- Nam, C., Bony, S., Dufresne, J., & Chepfer, H. (2012). The ‘too few, too bright’ tropical low-cloud problem in CMIP5 models. *Geophysical Research Letters*, 39(21), 2012GL053421. <https://doi.org/10.1029/2012GL053421>
- Ott, J., Pritchard, M., Best, N., Linstead, E., Curcic, M., & Baldi, P. (2020). *A Fortran-Keras deep learning bridge for scientific computing*. (Vol. 2020, pp. 1–13). Scientific Programming. <https://doi.org/10.1155/2020/8888811>
- Pahlavan, H. A., Hassanzadeh, P., & Alexander, M. J. (2024). Explainable offline-online training of neural networks for parameterizations: A 1D gravity wave-QBO testbed in the small-data regime. *Geophysical Research Letters*, 51(2), e2023GL106324. <https://doi.org/10.1029/2023GL106324>
- Pressel, K. G., Kaul, C. M., Schneider, T., Tan, Z., & Mishra, S. (2015). Large-eddy simulation in an elastic framework with closed water and entropy balances. *Journal of Advances in Modeling Earth Systems*, 7(3), 1425–1456. <https://doi.org/10.1002/2015MS000496>
- Ramadhan, A., Marshall, J., Souza, A., Lee, X. K., Piterbarg, U., Hillier, A., et al. (2023). Capturing missing physics in climate model parameterizations using neural differential equations. *arXiv*. (arXiv:2010.12559 [physics]). <https://doi.org/10.1002/essoar.10512533.1>
- Rasp, S. (2020). Coupled online learning as a way to tackle instabilities and biases in neural network parameterizations: General algorithms and Lorenz 96 case study (v1.0). *Geoscientific Model Development*, 13(5), 2185–2196. <https://doi.org/10.5194/gmd-13-2185-2020>
- Rasp, S., Pritchard, M. S., & Gentine, P. (2018). Deep learning to represent subgrid processes in climate models. *Proceedings of the National Academy of Sciences*, 115(39), 9684–9689. <https://doi.org/10.1073/pnas.1810286115>
- Romps, D. M. (2010). A direct measure of entrainment. *Journal of the Atmospheric Sciences*, 67(6), 1908–1927. <https://doi.org/10.1175/2010JAS3371.1>
- Savre, J. (2022). What controls local entrainment and detrainment rates in simulated shallow convection? *Journal of the Atmospheric Sciences*, 79(11), 3065–3082. <https://doi.org/10.1175/JAS-D-21-0341.1>
- Schillings, C., & Stuart, A. M. (2017). Analysis of the ensemble Kalman filter for inverse problems. *SIAM Journal on Numerical Analysis*, 55(3), 1264–1290. <https://doi.org/10.1137/16M105959X>

- Schneider, T., Leung, L. R., & Wills, R. C. J. (2024). Opinion: Optimizing climate models with process-knowledge, resolution, and AI. <https://doi.org/10.5194/egusphere-2024-20>
- Schneider, T., Stuart, A. M., & Wu, J.-L. (2021). Learning stochastic closures using ensemble Kalman inversion. *Transactions of Mathematics and Its Applications*, 5(1), tnab003. <https://doi.org/10.1093/imatrm/tnab003>
- Shamekh, S., & Gentine, P. (2023). Learning atmospheric boundary layer turbulence. <https://doi.org/10.22541/essoar.168748456.60017486/v1>
- Shankar, V., Puri, V., Balakrishnan, R., Maulik, R., & Viswanathan, V. (2023). Differentiable physics-enabled closure modeling for Burgers' turbulence. *Machine Learning: Science and Technology*, 4(1), 015017. <https://doi.org/10.1088/2632-2153/acb19c>
- Shen, C., Appling, A. P., Gentine, P., Bandai, T., Gupta, H., Tartakovsky, A., et al. (2023). Differentiable modelling to unify machine learning and physical models for geosciences. *Nature Reviews Earth and Environment*, 4(8), 552–567. <https://doi.org/10.1038/s43017-023-00450-9>
- Shen, Z. (2022). Data for "A library of large-eddy simulations forced by global climate models" (version 2.0) [Dataset]. *CaltechDATA*. <https://doi.org/10.22002/D1.20052>
- Shen, Z., Sridhar, A., Tan, Z., Jaruga, A., & Schneider, T. (2022). A library of large-eddy simulations forced by global climate models. *Journal of Advances in Modeling Earth Systems*, 14(3), e2021MS002631. <https://doi.org/10.1029/2021MS002631>
- Sherwood, S. C., Bony, S., & Dufresne, J.-L. (2014). Spread in model climate sensitivity traced to atmospheric convective mixing. *Nature*, 505(7481), 37–42. <https://doi.org/10.1038/nature12829>
- Siebesma, A. P., Soares, P. M. M., & Teixeira, J. (2007). A combined eddy-diffusivity/mass-flux approach for the convective boundary layer. *Journal of the Atmospheric Sciences*, 64(4), 1230–1248. <https://doi.org/10.1175/JAS3888.1>
- Siler, N., Po-Chedley, S., & Bretherton, C. S. (2018). Variability in modeled cloud feedback tied to differences in the climatological spatial pattern of clouds. *Climate Dynamics*, 50(3–4), 1209–1220. <https://doi.org/10.1007/s00382-017-3673-2>
- Soares, P., Miranda, P., Siebesma, A., & Teixeira, J. (2004). An eddy-diffusivity/mass-flux parametrization for dry and shallow cumulus convection. *Quarterly Journal of the Royal Meteorological Society*, 130(604), 3365–3383. <https://doi.org/10.1256/qj.03.223>
- Stevens, B., Lenschow, D. H., Vali, G., Gerber, H., Bandy, A., Blomquist, B., et al. (2003). Dynamics and chemistry of marine stratocumulus—DYCOMS-II. *Bulletin of the American Meteorological Society*, 84(5), 593. <https://doi.org/10.1175/BAMS-84-5-Stevens>
- Tan, Z., Kaul, C. M., Pressel, K. G., Cohen, Y., Schneider, T., & Teixeira, J. (2018). An extended eddy-diffusivity mass-flux scheme for unified representation of subgrid-scale turbulence and convection. *Journal of Advances in Modeling Earth Systems*, 10(3), 770–800. <https://doi.org/10.1002/2017MS001162>
- Thuburn, J., Weller, H., Vallis, G. K., Beare, R. J., & Whittall, M. (2018). A framework for convection and boundary layer parameterization derived from conditional filtering. *Journal of the Atmospheric Sciences*, 75(3), 965–981. <https://doi.org/10.1175/JAS-D-17-0130.1>
- Um, K., Brand, R., Fei, Y., Holl, P., & Thuerey, N. (2021). Solver-in-the-loop: Learning from differentiable physics to interact with iterative PDE-solvers. *Advances in Neural Information Processing Systems*, 33, 6111–6122.
- vanZanten, M. C., Stevens, B., Nuijens, L., Siebesma, A. P., Ackerman, A. S., Burnet, F., et al. (2011). Controls on precipitation and cloudiness in simulations of trade-wind cumulus as observed during RICO. *Journal of Advances in Modeling Earth Systems*, 3(2). <https://doi.org/10.1029/2011MS000056>
- Vial, J., Dufresne, J.-L., & Bony, S. (2013). On the interpretation of inter-model spread in CMIP5 climate sensitivity estimates. *Climate Dynamics*, 41(11–12), 3339–3362. <https://doi.org/10.1007/s00382-013-1725-9>
- Vignesh, P. P., Jiang, J. H., Kishore, P., Su, H., Smay, T., Brighton, N., & Velicogna, I. (2020). Assessment of CMIP6 cloud fraction and comparison with satellite observations. *Earth and Space Science*, 7(2), e2019EA000975. <https://doi.org/10.1029/2019EA000975>
- Wang, X., Han, Y., Xue, W., Yang, G., & Zhang, G. J. (2022). Stable climate simulations using a realistic general circulation model with neural network parameterizations for atmospheric moist physics and radiation processes. *Geoscientific Model Development*, 15(9), 3923–3940. <https://doi.org/10.5194/gmd-15-3923-2022>
- Watt-Meyer, O., Dresdner, G., McGibbon, J., Clark, S. K., Henn, B., Duncan, J., et al. (2023). Ace: A fast, skillful learned global atmospheric model for climate prediction. *arXiv*. (arXiv:2310.02074 [physics]).
- Webb, M. J., Andrews, T., Bodas-Salcedo, A., Bony, S., Bretherton, C. S., Chadwick, R., et al. (2017). The cloud feedback model intercomparison project (CFMIP) contribution to CMIP6. *Geoscientific Model Development*, 10(1), 359–384. <https://doi.org/10.5194/gmd-10-359-2017>
- Wu, J.-L., Levine, M. E., Schneider, T., & Stuart, A. (2023). Learning about structural errors in models of complex dynamical systems. *arXiv*. (arXiv:2401.00035 [physics]).
- Yu, S., Hannah, W., Peng, L., Lin, J., Bhouri, M. A., Gupta, R., et al. (2024). ClimSim: A large multi-scale dataset for hybrid physics-ML climate emulation. *arXiv*. (arXiv:2306.08754 [physics]).
- Yuval, J., & O'Gorman, P. A. (2020). Stable machine-learning parameterization of subgrid processes for climate modeling at a range of resolutions. *Nature Communications*, 11(1), 3295. <https://doi.org/10.1038/s41467-020-17142-3>
- Zelinka, M. D., Myers, T. A., McCoy, D. T., Po-Chedley, S., Caldwell, P. M., Ceppi, P., et al. (2020). Causes of higher climate sensitivity in CMIP6 models. *Geophysical Research Letters*, 47(1), e2019GL085782. <https://doi.org/10.1029/2019GL085782>
- Zhang, X.-L., Xiao, H., Luo, X., & He, G. (2022). Ensemble Kalman method for learning turbulence models from indirect observation data. *Journal of Fluid Mechanics*, 949, A26. <https://doi.org/10.1017/jfm.2022.744>