

1 A Fortran-Python Interface for Integrating Machine Learning Parameterization into
2 Earth System Models

3 Tao Zhang¹, Cyril Morcrette^{2,7}, Meng Zhang³, Wuyin Lin¹, Shaocheng Xie³, Ye Liu⁴, Kwinten Van
4 Weverberg^{5,6}, Joana Rodrigues²

- 5
6 1. Brookhaven National Laboratory, Upton, NY, USA
7 2. Met Office, FitzRoy Road, Exeter, EX13PB, UK
8 3. Lawrence Livermore National Laboratory, Livermore, CA, USA
9 4. Pacific Northwest National Laboratory, Richland, WA, USA
10 5. Department of Geography, Ghent University, Belgium
11 6. Royal Meteorological Institute of Belgium, Brussels, Belgium
12 7. Department of Mathematics and Statistics, Exeter University, Exeter, UK
13 Correspondence to: Tao Zhang (taozhang.ccs@gmail.com)

14 **Abstract**

15 Parameterizations in Earth System Models (ESMs) are subject to biases and uncertainties arising from
16 subjective empirical assumptions and incomplete understanding of the underlying physical processes.
17 Recently, the growing representational capability of machine learning (ML) in solving complex problems
18 has spawned immense interests in climate science applications. Specifically, ML-based parameterizations
19 have been developed to represent convection, radiation and microphysics processes in ESMs by learning
20 from observations or high-resolution simulations, which have the potential to improve the accuracies and
21 alleviate the uncertainties. Previous works have developed some surrogate models for these processes
22 using ML. These surrogate models need to be coupled with the dynamical core of ESMs to investigate
23 the effectiveness and their performance in a coupled system. In this study, we present a novel Fortran-
24 Python interface designed to seamlessly integrate ML parameterizations into ESMs. This interface
25 showcases high versatility by supporting popular ML frameworks like PyTorch, TensorFlow, and Scikit-
26 learn. We demonstrate the interface's modularity and reusability through two cases: a ML trigger function
27 for convection parameterization and a ML wildfire model. We conduct a comprehensive evaluation of
28 memory usage and computational overhead resulting from the integration of Python codes into the
29 Fortran ESMs. By leveraging this flexible interface, ML parameterizations can be effectively developed,
30 tested, and integrated into ESMs.

31

32 Plain Language

33 Earth System Models (ESMs) are crucial for understanding and predicting climate change. However, they
34 struggle to accurately simulate the climate due to uncertainties associated with parameterizing sub-grid
35 physics. Although higher-resolution models can reduce some uncertainties, they require significant
36 computational resources. Machine learning (ML) algorithms offer a solution by learning the important
37 relationships and features from high-resolution models. These ML algorithms can then be used to develop
38 parameterizations for coarser-resolution models, reducing computational and memory costs. To
39 incorporate ML parameterizations into ESMs, we develop a Fortran-Python interface that allows for
40 calling Python functions within Fortran-based ESMs. Through two case studies, this interface
41 demonstrates its feasibility, modularity and effectiveness.

42 1. Introduction

43 Earth System Models (ESMs) play a crucial role in understanding the mechanism of the climate system
44 and projecting future changes. However, uncertainties arising from parameterizations of sub-grid
45 processes pose challenges to the reliability of model simulations (Hourdin et al., 2017). Kilometer-scale
46 high-resolution models (Schär et al., 2020) can potentially mitigate the uncertainties by directly resolving
47 some key subgrid-scale processes that need to be parameterized in conventional low-resolution ESMs.
48 Another promising method, superparameterization – a type of multi-model framework (MMF) (D.
49 Randall et al., 2003; D. A. Randall, 2013), explicitly resolves sub-grid processes by embedding high-
50 resolution cloud-resolved models within the grid of low-resolution models. Consequently, both high-
51 resolution models and superparameterization approaches have shown promise in improving the
52 representation of cloud formation and precipitation. However, their implementation is challenged by
53 exceedingly high computational costs.

54
55 In recent years, machine learning (ML) techniques have emerged as a promising approach to
56 improve parameterizations in ESMs. They are capable of learning complex patterns and
57 relationships directly from observational data or high-resolution simulations, enabling the
58 capture of nonlinearities and intricate interactions that may be challenging to represent with
59 traditional parameterizations. For example, Zhang et al. (2021) proposed a ML trigger function
60 for a deep convection parameterization by learning from field observations, demonstrating its
61 superior accuracy compared to traditional CAPE-based trigger functions. Chen et al. (2023)
62 developed a neural network-based cloud fraction parameterization, better predicting both spatial

63 distribution and vertical structure of cloud fraction when compared to the traditional Xu-Randall
64 scheme (Xu & Randall, 1996). Krasnopolsky et al. (2013) prototyped a system using a neural
65 network to learn the convective temperature and moisture tendencies from cloud-resolving
66 model (CRM) simulations. These tendencies refer to the rates of change of various atmospheric
67 variables over one time step, diagnosed from particular parameterization schemes. These studies
68 lay the groundwork for integrating ML-based parameterization into ESMs.

69
70 However, the aforementioned studies primarily focus on offline ML of parameterizations that do
71 not directly interact with ESMs. Recently, there have been efforts to implement ML
72 parameterizations that can be directly coupled with ESMs. Several studies have developed ML
73 parameterizations in ESMs by hard coding custom neural network modules, such as O’Gorman
74 & Dwyer (2018), Rasp et al. (2018), Han et al. (2020) and Gettelman et al. (2021). They
75 incorporated a Fortran-based ML inference module to allow the loading of the pre-trained ML
76 weights to reconstruct the ML algorithm in ESMs. The hard-coding has limitations. [When a
77 trained ML model is incorporated into ESMs, it is frozen and cannot be updated during runtime.
78 Recently, Kochkov et al.\(2024\) introduced the NeuralGCM, an innovative approach that enables
79 the ML model to be updated during runtime with a differentiable dynamical core. This allows for
80 end-to-end training and optimization of the interactions with large-scale dynamics. However, the
81 hard-coding coupling method does not support such capability.](#)

82
83 Fortran-Keras Bridge (FKB; Ott et al. (2020)) and C Foreign Function Interface (CFFI;
84 <https://cffi.readthedocs.io>) are two packages that support two-way coupling between Fortran-based ESM
85 and Python based ML parameterizations. FKB enables tight integration of Keras deep learning models but
86 is specifically bound to the Keras library, limiting its compatibility with other frameworks like PyTorch
87 and Scikit-Learn. On the other hand, CFFI provides a more flexible solution that in principle supports
88 coupling various ML packages due to its language-agnostic design. Brenowitz & Bretherton (2018)
89 utilized it to enable the calling of Python ML algorithms within ESMs. However, the CFFI has several
90 limitations. When utilizing CFFI to interface Fortran and Python, it uses global data structures to pass
91 variables between the two languages. This approach results in additional memory overhead as variable
92 values need to be copied between languages, instead of being passed by reference. Additionally, CFFI
93 lacks automatic garbage collection for the unused memory within these data structures and copies.
94 Consequently, the memory usage of the program gradually increases over its lifetime. In addition, when

Deleted: (

Deleted: ,

Deleted: S

Deleted: uch hard-coding approach restricts the ML algorithm’s ability to adapt to changes in the model dynamics over time, as the ‘online’ updating requires a two-way coupling between the dominantly Fortran-based ESMs and Python ML libraries.

103 using CFFI to call Python functions from a Fortran program, the process involves several steps such as
104 registering variables into a global data structure, calling the Python function, and retrieving the calculated
105 result. These multiple steps can introduce computational overhead due to the additional operations
106 required.

107
108 Additionally, Wang et al. (2022) developed a coupler to facilitate two-way communication between ML
109 parameterizations and host ESMs. The coupler gathers state variables from the ESM using the Message
110 Passing Interface (MPI) and transfers them to a Python-based ML module. It then receives the output
111 from the Python code and returns them to the ESM. While this approach effectively bridges Fortran and
112 Python, its use of file-based data passing to exchange information between modules carries some
113 performance overhead relative to tighter coupling techniques. Optimizing the data transfer, such as via
114 shared memory, remains an area for improvement to fully leverage this coupler's ability to integrate
115 online-adaptive ML parameterizations within large-scale ESM simulations, which is the main goal for this
116 study.

117
118 In this study, we investigate the integration of ML parameterizations into Fortran-based ESM
119 models by establishing a flexible interface that enables the invocation of ML algorithms in
120 Python from Fortran. This integration offers access to any Python codes from Fortran, including
121 a diverse range of ML frameworks, such as PyTorch, TensorFlow, and Scikit-learn, which can
122 effectively be utilized for parameterizing intricate atmospheric and other climate system
123 processes. The coupling of the Fortran model and the Python ML code needs to be performed for
124 thousands of model columns and over thousands of timesteps for a typical model simulation.
125 Therefore, it is crucial for the coupling interface to be both robust and efficient. We showcase the
126 feasibility and benefits of this approach through case studies that involve the parameterization of
127 deep convection and wildfire processes in ESMs. The two cases demonstrate the robustness and
128 efficiency of the coupling interface. The focus of this paper is on documenting the coupling
129 between the Fortran ESM and the ML algorithms and systematically evaluating the
130 computational efficiency and memory usage of different ML frameworks (such as Pytorch and
131 TensorFlow), different ML algorithms, and different configuration of a climate model. The
132 assessment of the scientific performance of the ML emulators will be addressed in follow-on
133 papers. The showcase examples emphasize the potential for high modularity and reusability by
134 separating the ML components into Python modules. This modular design facilitates independent

135 development and testing of ML-based parameterizations by researchers. It enables easier code
136 maintenance, updates, and the adoption of state-of-the-art ML techniques with only minimal
137 disruption to the existing Fortran infrastructure. Ultimately, this advancement will contribute to
138 enhanced predictions and a deeper comprehension of the evolving climate of our planet. It is
139 important to note that the current interface only supports executing deep learning algorithms on CPUs and
140 does not support running them on GPUs.

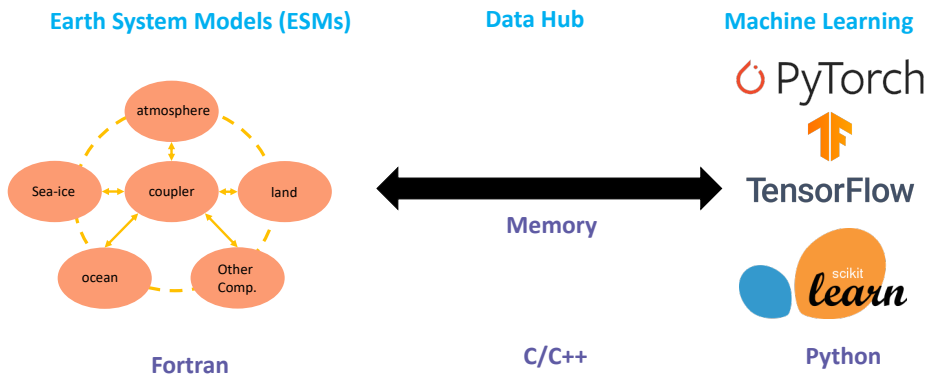
Deleted: ing

141
142 The rest of this manuscript is organized as follows: Section 2 presents the detailed interface that
143 integrates ML into Fortran-based ESM models. Section 3 discusses the performance of the
144 interface and presents its application in two case studies. Finally, Section 4 provides a summary
145 of the findings and a discussion of their implications.

146 2. General design of the ML interface

147 2.1 Architecture of the ML interface

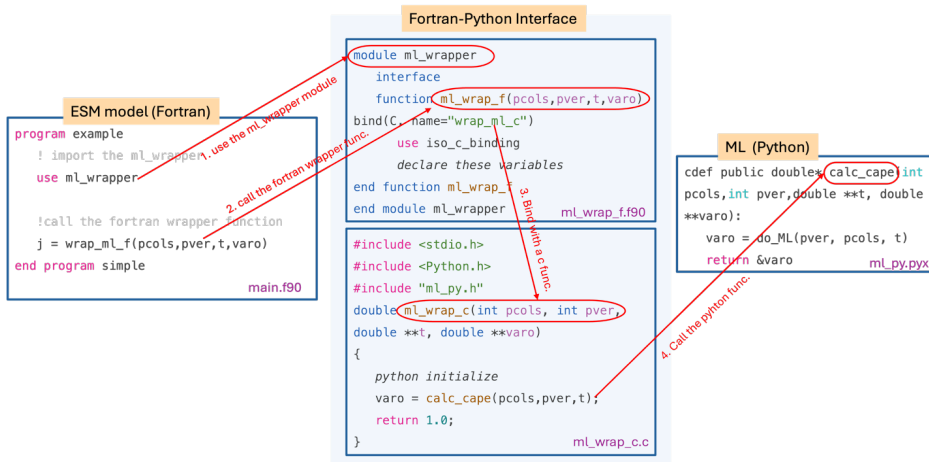
148 We developed an interface using shared memory to enable two-way coupling between Fortran and Python
149 (Figure 1). The ESM used in the demonstration in Figure 1 is the U.S. Department of Energy (DOE)
150 Energy Exascale Earth System Model (E3SM; Golaz et al., 2019, 2022). Because Fortran cannot directly
151 call Python, we utilized C as an intermediary since Fortran can call C functions. This approach leverages
152 C as a data hub to exchange information without requiring a framework-specific binding like KFB. As a
153 result, our interface supports invoking any Python-based ML package such as PyTorch, TensorFlow, and
154 scikit-learn from Fortran. While C can access Python scalar values through the built-in
155 PyObject_CallObject function from the Python C API, we employed Cython for its ability to transfer
156 array data between the languages. Using Cython, multidimensional data structures can be efficiently
157 passed between Fortran and Python modules via C, allowing for flexible training of ML algorithms within
158 ESMs.



160
 161 **Figure 1.** The interface of the ML bridge for two-way communication via memory between Fortran ESM
 162 and Python ML module.

163 2.2 Code structure

164 Figure 2 illustrates how the framework operates using a toy code example. The Fortran-Python interface
 165 comprises a Fortran wrapper and C wrapper files, which are bound together. The Fortran-based ESM first
 166 imports the Fortran wrapper, allowing it to call wrapper functions with input and output memory
 167 addresses. The interface then passes these memory addresses to the Python-based ML module, which
 168 performs the ML predictions and returns the output address to the Fortran model.

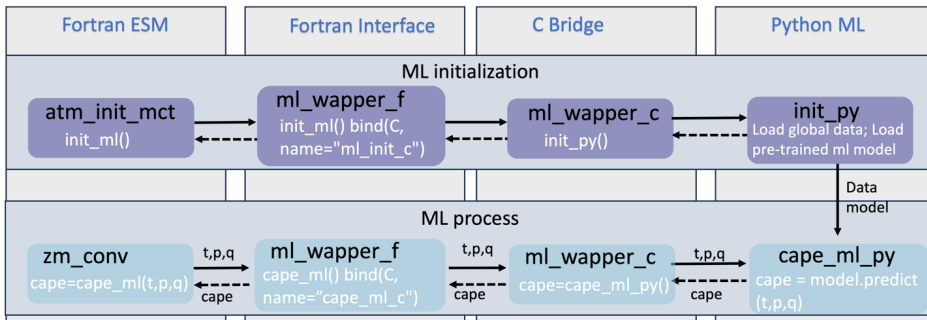


169

170 Figure 2. Toy code illustrating the Fortran-Python interface. [It is noted that a fleshed-out, com-
171 pliable version of this toy example exists in the linked GitHub repository.](#)
172

173 When coupling the Python ML module with the [Fortran](#) model using the interface, additional steps should
174 be considered: 1. The ML module should remain active throughout the model simulations, without any
175 Python finalization calls, ensuring it is continuously available. 2. The Python module should load the
176 trained ML model and any required global data only once, rather than at each simulation step. This one-
177 time initialization process improves efficiency and prevents unnecessary repetition. On the Fortran ESM
178 side, the `init_ml()` function is called within the `atm_init_mct` module to load the ML model and global
179 data (shown in Figure 3). Then, similar to the toy code, we call the wrapper function, pass input variables
180 to Python for ML predictions, and return the results to the Fortran side. 3. When compiling the complex
181 system, which includes Python, C, Cython, and Fortran code, the Python path should be specified in the
182 CFLAGS and LDFLAGS. It is important to note that without the position-independent compiling flag (`-
183 fPIC`), the hybrid system will only work on a single node and may cause segmentation faults on multiple
184 nodes. Including it can resolve this issue, allowing multi-node compatibility.
185
186

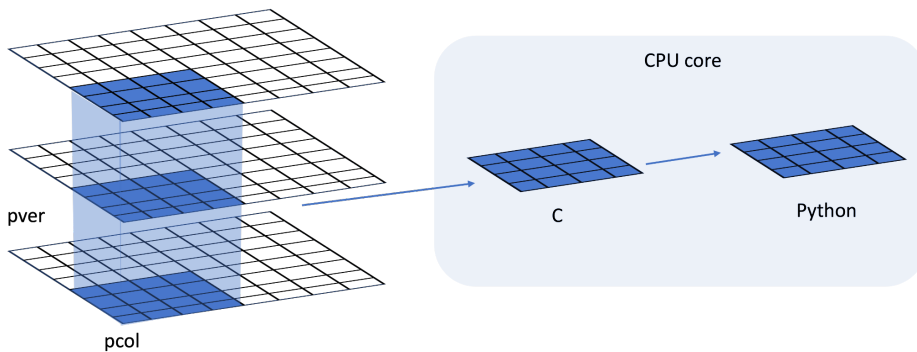
Deleted: real



187 **Figure 3.** The code structure of the ML bridge interface using the ML closure in deep convection as an
188 example.
189

190
191 In traditional ESMs, sub-grid scale parameterization routines such as convection parameterizations are
192 often calculated separately for each vertical column of the model domain. Meanwhile, the domain is
193 typically decomposed horizontally into 2D chunks that can be solved in parallel using MPI processes.
194 Each CPU core/MPI process is assigned a number of chunks of model columns to update asynchronously
195 (Figure 4). Our interface takes advantage of this existing parallel decomposition by designing the ML

197 calls to operate over all columns simultaneously within each chunk, rather than invoking the ML scheme
198 individually for each column. This allows the coupled model-ML system to leverage parallelism in the
199 neural network computations. If the ML were called separately for every column, parallel efficiencies
200 would not be realized. By aggregating inputs over the chunk-scale prior to interfacing with Python,
201 performance is improved through better utilization of multi-core and GPU-based ML capabilities during
202 parameterization calculations.
203



204
205 **Figure 4.** Data and system structure. The model domain is decomposed into chunks of columns. pver
206 refers to number of pressure vertical levels. A chunk contains multiple columns (up to pcol). Multiple
207 chunks can be assigned to each CPU core.
208
209

210 3. Results

211 The framework explained in the previous section provides seamless support for various ML
212 parameterizations and various ML frameworks, such as PyTorch, Tensorflow, and Scikit-learn. To
213 demonstrate the versatility of this framework, we applied it in two distinct case applications. The first
214 application replaces the conventional CAPE-based trigger function in a deep convection parameterization
215 with a machine-learned trigger function. The second application involves a ML-based wildfire model that
216 interacts bidirectionally with the ESM. We provide a brief introduction to these two cases. Detailed
217 descriptions and evaluations will be presented in separate papers.
218

219 The framework's performance is influenced by two primary factors: increasing memory usage and
220 increasing computational overhead. Firstly, maintaining the Python environment fully persistent in
221 memory throughout model simulations can impact memory usage, especially for large ML algorithms.
222 This elevated memory footprint increases the risk of leaks or crashes as simulations progress. Secondly,
223 executing ML components within the Python interpreter inevitably introduces some overhead compared
224 to the original ESMs. The increased memory requirements and decreased computational efficiency
225 associated with these considerations can impact the framework's usability, flexibility, and scalability for
226 different applications.

227

228 To comprehensively assess performance, we conducted a systematic evaluation of various ML
229 frameworks, ML algorithms, and physical models. This evaluation is built upon the foundations
230 established for evaluating the ML trigger function in the deep convection parameterization.

231 3.1 Application cases

232 3.1.1 ML trigger function in deep convection parameterization

233 In General Circulation Models, uncertainties in convection parameterizations are recognized to be closely
234 linked to the convection trigger function used in these schemes (Bechtold et al., 2004; Xie et al., 2004,
235 2019; Xie & Zhang, 2000; Lee et al., 2007). The convective trigger in a convective parameterization
236 determines when and where model convection should be triggered as the simulation advances. In many
237 convection parameterizations, the trigger function consists of a simple, arbitrary threshold for a physical
238 quantity, such as convective available potential energy (CAPE). Convection will be triggered if the CAPE
239 value exceeds a threshold value.

240

241 In this work, we use this interface to test a newly developed ML trigger function in E3SM. The ML
242 trigger function was developed with the training data originating from simulations performed using the
243 kilometer-resolution (1.5 km grid spacing), Met Office Unified Model Regional Atmosphere 1.0
244 configuration (Bush et al., 2020). Each simulation consists of a limited area model (LAM) nested within a
245 global forecast model providing boundary conditions (Walters et al., 2017; Webster et al., 2008). In total
246 80 LAM simulations were run located so as to sample different geographical regions worldwide. Each
247 LAM was run for 1 month, with 2-hourly output, using a grid-length of 1.5 km, a 512 x 512 domain, and
248 a model physics package used for operational weather forecasting. The 1.5 km data is coarse-grained to
249 several scales from 15 to 144 km.

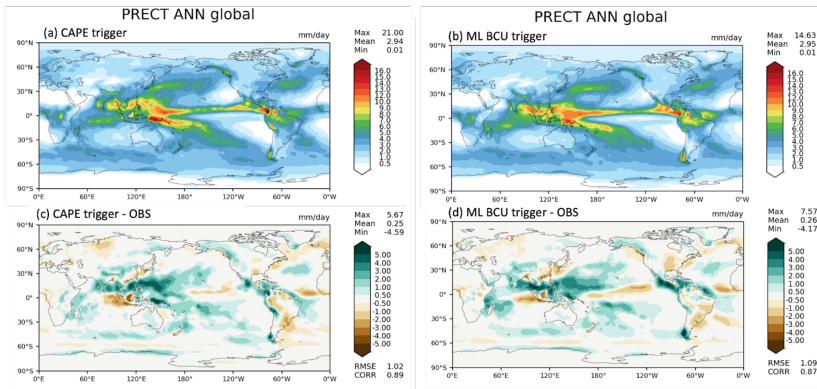
250

Deleted: .

252 A two-stream neural network architecture is used for the ML model. The first stream takes profiles of
 253 temperature, specific humidity and pressure across 72 levels at each scale as inputs and passes them
 254 through a 4-layer convolutional neural network (CNN) with kernel sizes of 3, to extract large scale
 255 features. The second stream takes mean orographic height, standard deviation of orographic height, land
 256 fraction and the size of the grid-box as inputs. The outputs of the two streams are then combined and fed
 257 into a 2-layer fully connected network to allow the ML model to leverage both atmospheric and surface
 258 features when making its predictions. The output is a binary variable indicating whether the convection
 259 happens, based on the condition of buoyant cloudy updrafts (BCU, e.g. Hartmann et al., 2019; Swann,
 260 2001). If there are 3 contiguous levels where the predicted BCU is larger than 0.05, the convection
 261 scheme is triggered. Once trained, the CNN is coupled to E3SM and thermodynamic information from
 262 E3SM is passed to it to predict the trigger condition. Then, the predicted result is returned to E3SM.

263
 264 Figure 5 shows the comparison of annual mean precipitation between the control run using the traditional
 265 CAPE-based trigger function and the run using the ML BCU trigger function. The ML BCU scheme
 266 demonstrates reasonable spatial patterns of precipitation, similar to the control run, with comparable root-
 267 mean-square error and spatial correlation. Additional experiments exploring the definition of BCU and
 268 varying the thresholds along with an in-depth analysis will be presented in a follow-up paper.

269



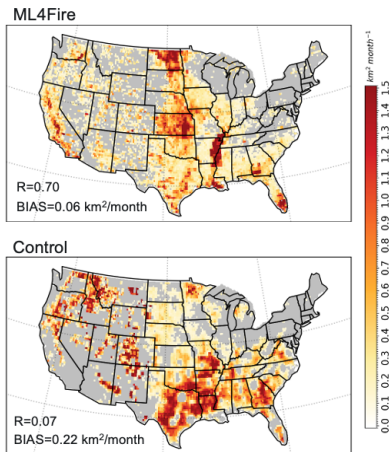
270
 271 **Figure 5.** Comparison of annual mean precipitation between the control run using the CAPE-based
 272 trigger function (a, c) and the run using the ML BCU trigger function (b, d).

273 3.1.2 ML learning fire model

274 Predicting wildfire burned area is challenging due to the complex interrelationships between fires,
275 climate, weather, vegetation, topography, and human activities (Huang et al., 2020). Traditionally,
276 statistical methods like multiple linear regression have been applied, but are limited in the number and
277 diversity of predictors considered (Yue et al., 2013). In this study, we develop a coupled fire-land-
278 atmosphere framework that uses machine learning to predict wildfire area, enhancing long-term burned
279 area projections and assessing fire impacts by enabling simulations of interactions among fire,
280 atmosphere, land cover, and vegetation.

281
282 The ML algorithm is trained using a monthly dataset, which includes the target variable of burned area, as
283 well as various predictor variables. These predictors encompass local meteorological data (e.g., surface
284 temperature, precipitation), land surface properties (e.g., monthly mean evapotranspiration and surface
285 soil moisture), and socioeconomic variables (e.g., gross domestic product, population density), as
286 described by Wang et al. (2022). In the coupled fire-land-atmosphere framework, meteorology variables
287 and land surface properties are provided by the E3SM. We use the eXtreme Gradient Boosting algorithm
288 implemented in Scikit-Learn to train the ML fire model. Figure 6 demonstrates that the ML4Fire model
289 exhibits superior performance in terms of spatial distribution compared to process-based fire models,
290 particularly in the Southern US region. Detailed analysis will be presented in a separate paper. The
291 ML4Fire model has proven to be a valuable tool for studying vegetation-fire interactions, enabling
292 seamless exploration of climate-fire feedbacks.

293
294

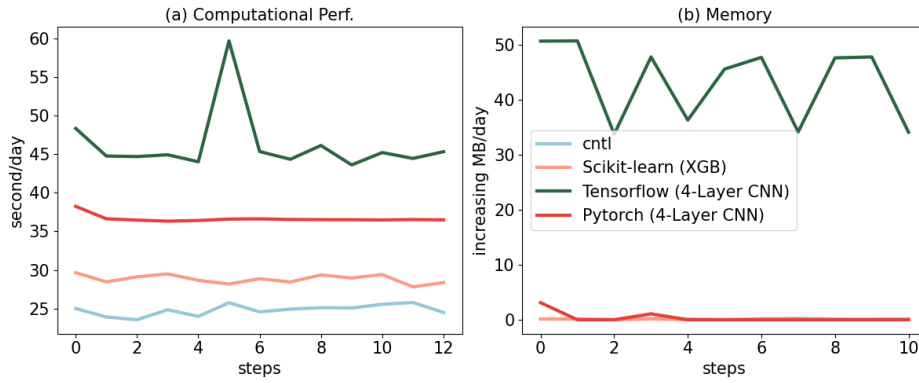


295

296 **Figure 6.** Comparison between ML4Fire model and process-based fire model against the historical
 297 burned area from Global Fire Emissions Database 5 from 2001-2020. R and BIAS are the spatial
 298 pattern correlation and difference against the observation, respectively.

299 **3.2 Performance of different ML frameworks**

300 The Fortran-Python bridge ML interface supports various ML frameworks, including PyTorch,
 301 TensorFlow, and scikit-learn. These ML frameworks can be trained offline using kilometer-scale high-
 302 resolution models (such as the ML trigger function) or observations (ML fire model). Once trained, they
 303 can be plugged into the ML bridge interface through different API interfaces specific to each framework.
 304 The coupled ML algorithms are persistently resident in memory, just like the other ESM components.
 305 During each step of the process, the performance of the full system is significantly affected by memory
 306 usage. If memory consumption increases substantially, it may lead to memory leaks as the number of time
 307 step iteration increases. In addition, Python, being an interpreted language, is typically considered to have
 308 slower performance compared to compiled languages like C/C++ and Fortran. Therefore, incorporating
 309 Python may decrease computational performance. We examine the memory usage and computational
 310 performance across various ML frameworks based on implementing the ML trigger function in E3SM.
 311 The ML algorithm is implemented as a two-stream CNN model using Pytorch and TensorFlow
 312 frameworks, as well as XGBoost using the Scikit-learn package. It should be noted that XGBoost, a
 313 boosting tree-based model, is a completely different type of ML model compared to the CNNs, which are
 314 the type of deep neural network.



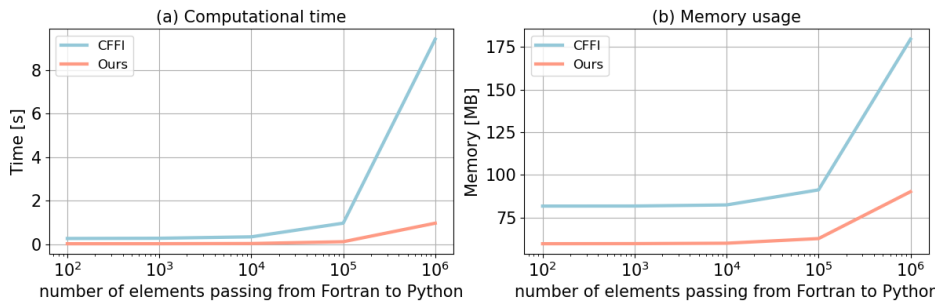
315
 316 **Figure 7.** Computational and memory overhead as the simulation progresses for coupling the ML trigger
 317 function with the E3SM model. The x-axis represents the simulated time step. The y-axis of (a) represents
 318 the simulation speed measured in seconds per day (indicating the number of seconds required to simulate
 319 one day). The y-axis of (b) represents the relative increase in memory usage for Scikit-learn, TensorFlow,
 320 and PyTorch compared with CNTL. CNTL represents the original simulation without using the ML
 321 framework.

322
 323 Figure 7 illustrates the computational and memory overhead associated with the ML parameterization
 324 using different ML frameworks. It shows that XGBoost only exhibits a 20% increase in the simulation
 325 time required for simulating one day due to its simpler algorithm. For more complex neural networks,
 326 PyTorch incurs a 52% overhead, while TensorFlow's overhead is almost 100% – about two times as much
 327 as the overhead by PyTorch. In terms of memory usage, we use the highwater memory metric (Gerber &
 328 Wasserman, 2013), which represents the total memory footprint of a process. Scikit-learn and PyTorch do
 329 not show any significant increase in memory usage. However, TensorFlow shows a considerable increase
 330 up to 50MB per simulation day per MPI process element. This is significant because for a node with 48
 331 cores, it would equate to an increase of around 2GB per simulated day on that node. This rapid memory
 332 growth could quickly lead to a simulation crash due to insufficient memory during continuous
 333 integrations, preventing the use in practical simulations. Our findings show that the TensorFlow
 334 prediction function does not release memory after each call. Therefore, we recommend using PyTorch for
 335 complex deep learning algorithms and Scikit-learn for simpler ML algorithms to avoid these potential
 336 memory-related issues when using TensorFlow.

337
 338 Previous work, such as Brenowitz & Bretherton (2018, 2019) has utilized the CFFI package to establish
 339 communication between Fortran ESM and ML Python. As described in the Introduction, while CFFI
 340 offers flexibility in supporting various ML packages, it does have certain limitations. To pass variables

341 from Fortran to Python, the approach relies on global data structures to store all variables, including both
 342 the input from Fortran to Python and the output returning to Fortran. Consequently, this package results in
 343 additional memory copy operations and increasing overall memory usage. In contrast, our interface takes
 344 a different approach by utilizing memory references to transfer data between Fortran and Python,
 345 avoiding the need for global data structures and the associated overhead. This allows for a more efficient
 346 data transfer process.

347
 348 In Figure 8, we present a comparison between the two frameworks by testing the different number of
 349 elements passed from Fortran to Python. The evaluation is based on a demo example that focuses solely
 350 on declaring arrays and transferring them from Fortran to Python, rather than a real E3SM simulation.
 351 Figure 8a illustrates the impact of the number of passing elements on the overhead of the two interfaces.
 352 As the number of elements exceeds 10^4 , the overhead of CFFI becomes significant. When the number
 353 surpasses 10^6 , the overhead of CFFI is nearly ten times greater than that of our interface. Regarding
 354 memory usage, our interface maintains a stable memory footprint of approximately 60MB. Even as the
 355 number of elements increases, the memory usage only shows minimal growth. However, for CFFI, the
 356 memory usage starts at 80MB, which is 33% higher than our interface. As the number of elements
 357 reaches 10^6 , the memory overhead for CFFI dramatically rises to 180MB, twice as much as our interface.
 358



359
 360 Figure 8. Comparison of our framework and the CFFI framework in terms of computational time
 361 and memory usage. The x-axis represents the number of elements transferred from Fortran to
 362 Python, while the y-axis displays the total time (a) and total memory usage (b) for a
 363 demonstration example. The evaluations presented are based on the average results obtained
 364 from 5 separate tests.
 365

366 3.3 Performance of ML algorithms of different complexities

367 ML parameterizations can be implemented using various deep learning algorithms with different levels of
368 complexity. The computational performance and memory usage can be influenced by the complexity of
369 these algorithms. In the case of the ML trigger function, a two-stream four-layer CNN structure is
370 employed. We compare this structure with other ML algorithms such as Artificial Neural Network (ANN)
371 and Residual Network (ResNet), whose structures are detailed in Table 1. We selected these three ML
372 algorithms because they are commonly used in previous ML parameterization approaches, such as
373 (Brenowitz & Bretherton, 2019; Han et al., 2020; Wang et al., 2022). Systematically evaluating the hybrid
374 system with these ML methods using our interface can help identify bottlenecks and improve the system
375 computational performance. These algorithms are implemented in PyTorch. The algorithm’s complexity
376 is measured by the number of parameters, with the CNN having approximately 60 times more parameters
377 than ANN, and ResNet having roughly 1.5 times more parameters than CNN.

378
379
380

381 **Table 1.** The structure and number of parameters of each ML algorithms.

Algorithms	Structure	# of parameters
ANN	3 x Linear	121,601
CNN	4 x Conv2d + 2 x Linear	7,466,753
ResNet	17 x Conv2d + 1 x Linear	11,177,025

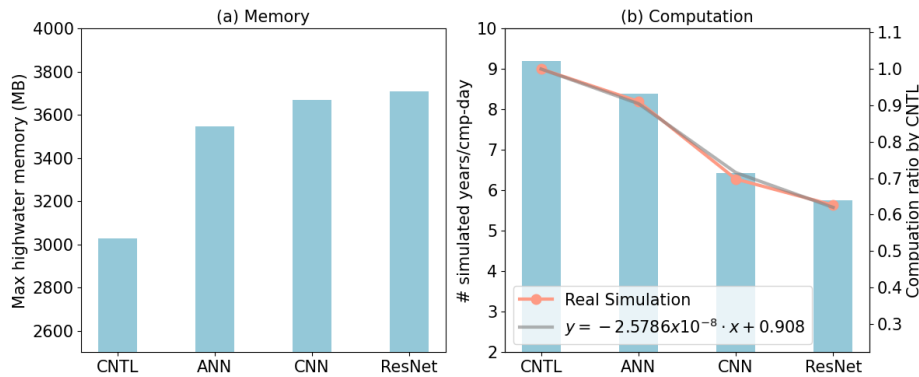
382
383 Figure 9 presents a comparison of the memory and computational costs between the CNTL run without
384 deep learning parameterization and the hybrid run with various deep learning algorithms. The same
385 specific process-element layout (placement of ESM component models on distributed CPU cores) is used
386 for all the simulations. Deep learning algorithms incur a significant yet affordable increase in memory
387 overhead, with at least a 20% increase compared to the CNTL run (Figure 9a). This is primarily due to the
388 integration of ML algorithms into the ESM, which persist throughout the simulations. Although there is a
389 notable increase in complexity among the deep learning algorithms, their memory usage only shows a
390 slight rise. This is because the memory increment resulting from the ML parameters is relatively small.
391 Specifically, ANN requires 1MB of memory, CNN requires 60MB, and the ResNet algorithms requires
392 85MB, which are calculated based on the number of parameters in each algorithm. When comparing these

393 values to the memory consumption of the CNTL run, which is approximately 3000MB, the additional
 394 parameters' incremental memory consumption is not substantial. However, when we use 128 MPI
 395 processes per node, it could bring the total memory requirement to approximately 460 GB per node. If the
 396 available hardware memory is less than this, the process layout must be adjusted accordingly.

397
 398 In terms of computational performance, the Python-based ML calls inevitably introduce some overhead.
 399 However, as shown in Figure 9b, the performance decrease is not substantial. The simple ANN model
 400 reduces performance by only about 10% compared to the CNTL run, while even the more complex
 401 ResNet model results in a 35% decrease. In contrast, Wang et al. (2022) reported a 100% overhead in
 402 their interface [when using the ResNet model as well](#), which transfers parameters via files. It is worth
 403 noting that in this study, the deep learning algorithms are executed on CPUs. To enhance computational
 404 performance, future work could consider utilizing GPUs for acceleration.

405
 406 In addition, we develop a performance model to estimate computational performance for the hybrid
 407 model using different ML model sizes and complexities. This performance model, based on linear
 408 regression, [predicts the ratio of the simulated years per day of the ML-augmented run to that of the CNTL](#)
 409 [run as a function of the number of ML parameters](#), shown in Figure 9b. It provides a simple yet effective
 410 way to capture this relationship and serves as a valuable tool for performance prediction when
 411 incorporating more complicated ML models.

Deleted: predicts the computational ratio relative to the CNTL run by taking the number of ML parameters as input

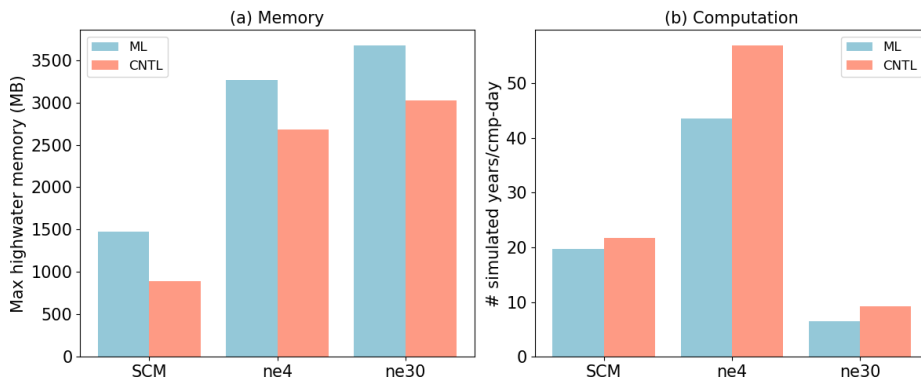


412
 413 **Figure 9.** Comparison of CNTL and the hybrid model using various ML algorithms in terms of memory
 414 and computation. CNTL is the default run without ML parameterizations. In (b), the left y-axis represents
 415 the actual number of simulated years per day, while the right y-axis shows the relative performance

418 compared to the CNTL run (orange line). The gray line illustrates the regression between the number of
 419 ML parameters (x) and the relative performance of the hybrid system (y).
 420

421 3.4 Performance for physical models of different complexities

422 ML parameterization can be applied to various ESM configurations, for example, with the E3SM
 423 Atmosphere Model (EAM), we experiment with Single Column Model (SCM), the ultra low-resolution
 424 model of EAM (ne4), and the nominal low resolution model of EAM (ne30) configurations. The SCM
 425 consists of one single atmosphere column of a global EAM (Bogenschutz et al., 2020; Gettelman et al.,
 426 2019). ne4 has 384 columns, with each column representing the horizontal resolution of 7.5°. ne30 is the
 427 default resolution for EAM and comprises 21,600 columns, with each column representing the horizontal
 428 resolution of 1°. In the case of the ML trigger function, the memory overhead is approximately 500MB
 429 for all configurations due to the loading of the ML algorithm, which does not vary with the configuration
 430 of the ESM.
 431



432
 433 **Figure 10.** Comparison of CNTL and ML for various ESMs in terms of memory and computation. The
 434 ESM configuration include SCM, ultra-low resolution model (ne4) and nominal low-resolution model
 435 (ne30).

436
 437 Regarding computational performance, SCM utilizes 1 process, ne4 employs 1 node with 64 processes,
 438 and ne30 utilizes 10 nodes with each node using 128 processes. In the case of SCM, the overhead
 439 attributed to the ML parameterization is approximately 9% due to the utilization of only 1 process.
 440 However, for ne4 and ne30, the overhead is 23% and 28% respectively (Figure 10). The increasing

Deleted: ss

442 computational overhead is primarily due to resource competition when multiple processes are used within
443 a single node. It is noted that although there is a significant computational gap between ML and CNTL
444 for ne4, the relative performance between ML and CNTL for ne4 is approximately 76.7%, which is close
445 to ne30 at 71.4%.

446

447 4. Discussion and Conclusion

448 ML algorithm can learn detailed information about cloud processes and atmospheric dynamics from
449 kilometer-scale models and observations and serves as an approximate surrogate for the kilometer-scale
450 model. Instead of explicitly simulating kilometer-scale processes, the ML algorithms can be designed to
451 capture the essential features and relationships between atmospheric variables by training on available
452 kilometer-scale data. The trained algorithms can then be used to develop parameterizations for use in
453 models at coarser resolutions, reducing the computational and memory costs. By using ML
454 parameterizations, scientists can effectively incorporate the insights gained from kilometer-scale models
455 for coarser-resolution simulations. Through learning the complex relationships and patterns present in the
456 high-resolution data, the ML-based parameterizations have the potentials to more accurately represent
457 cloud processes and atmospheric dynamics in the ESMs. This approach strikes a balance between
458 computational efficiency and capturing critical processes, enabling more realistic simulations and
459 predictions while minimizing computational resources. All these potential benefits in turn promote
460 innovative developments to facilitate increasing and more efficient use of ML parameterizations.

461

462 In this study, we develop a novel Fortran-Python interface for developing ML parameterizations. This
463 interface demonstrates feasibility in supporting various ML frameworks, such as PyTorch, TensorFlow,
464 and Scikit-learn and enables the effective development of new ML-based parameterizations to explore
465 ML-based applications in ESMs. Through two cases - a ML trigger function in convection
466 parameterization and a ML wildfire model - we highlight high modularity and reusability of the
467 framework. We conduct a systematic evaluation of memory usage and computational overhead from the
468 integrated Python codes.

469

470 Based on our performance evaluation, we observe that coupling ML algorithms using TensorFlow into
471 ESMs can lead to memory leaks. As a recommendation, we suggest using PyTorch for complex deep
472 learning algorithms and Scikit-learn for simple ML algorithms for the Fortran-Python ML interface.

473

474 The memory overhead primarily arises from loading ML algorithms into ESMs. If the ML algorithms are
475 implemented using PyTorch or Scikit-learn, the memory usage will not increase significantly. The
476 computational overhead is influenced by the complexity of the neural network and the number of
477 processes running on a single node. As the complexity of the neural network increases, more parameters
478 in the neural network require forward computation. Similarly, when there are more processes running
479 on a single node, the integrated Python codes introduce more resource competition.

480
481 Although this interface provides a flexible tool for ML parameterizations, it does not currently utilize
482 GPUs for ML algorithms. In Figure 3, it is shown that each chunk is assigned to a CPU core. However, to
483 effectively leverage GPUs, it is necessary to gather the variables from multiple chunks and pass them to
484 the GPUs. Additionally, if an ESM calls the Python ML module multiple times in each time step, the
485 computational overhead becomes significant. It is crucial to gather the variables and minimize the number
486 of calls. In the future, we will enhance the framework to support this mechanism, enabling GPU
487 utilization and overall performance improvement.

488 Acknowledge

489 This work was primarily supported by the Energy Exascale Earth System Model (E3SM) project of the
490 Earth and Environmental System Modeling program, funded by the US Department of Energy, Office of
491 Science, Office of Biological and Environmental Research. Research activity at BNL was under the
492 Brookhaven National Laboratory contract DE-SC0012704 (Tao Zhang, Wuyin Lin). The work at LLNL
493 was performed under the auspices of the US Department of Energy by the Lawrence Livermore National
494 Laboratory under Contract DE-AC52-07NA27344. The work at PNNL is performed under the Laboratory
495 Directed Research and Development Program at the Pacific Northwest National Laboratory. PNNL is
496 operated by DOE by the Battelle Memorial Institute under contract DE-A05-76RL01830.

497

498 Author contribution

499 TZ developed the Fortran-Python Interface. CM and JR contributed the ML model for the trigger
500 function. YL contributed the ML model for the wire fire model. TZ and MZ assessed the performance of
501 the ML trigger function. TZ took the lead in preparing the manuscript, with valuable edits from CM, MZ,
502 WL, SX, YL, KW, and JR. All the co-authors provided valuable insights and comments for the
503 manuscript.

504 **Conflict of Interest**

505 The authors declare that they have no conflict of interest.

506

507 **Data Availability Statement**

508 The Fortran-Python interface for developing ML parameterizations can be archived at
509 <https://doi.org/10.5281/zenodo.11005103> (Zhang et al., 2024) and can be also accessed at
510 <https://github.com/tzhang-ccs/ML4ESM>. The E3SM model can be accessed at
511 <https://zenodo.org/records/12175988>. The dataset for machine learning trigger function can be accessed
512 at <https://zenodo.org/records/12205917>. The dataset for machine learning wild fire can be accessed at
513 <https://zenodo.org/records/12212258>.

514 **References**

515 Bechtold, P., Chaboureau, J.-P., Beljaars, A., Betts, A. K., Köhler, M., Miller, M., & Redelsperger, J.-L.

516 (2004). The simulation of the diurnal cycle of convective precipitation over land in a global

517 model. *Quarterly Journal of the Royal Meteorological Society*, 130(604), 3119–3137.

518 <https://doi.org/10.1256/qj.03.103>

519 Bogenschutz, P. A., Tang, S., Caldwell, P. M., Xie, S., Lin, W., & Chen, Y.-S. (2020). The E3SM version

520 1 single-column model. *Geoscientific Model Development*, 13(9), 4443–4458.

521 <https://doi.org/10.5194/gmd-13-4443-2020>

522 Brenowitz, N. D., & Bretherton, C. S. (2018). Prognostic validation of a neural network unified physics

523 parameterization. *Geophysical Research Letters*, 45(12), 6289–6298.

524 <https://doi.org/10.1029/2018gl078510>

525 Brenowitz, N. D., & Bretherton, C. S. (2019). Spatially extended tests of a neural network

526 parametrization trained by coarse-graining. *Journal of Advances in Modeling Earth Systems*,

527 11(8), 2728–2744. <https://doi.org/10.1029/2019ms001711>

528 Bush, M., Allen, T., Bain, C., Boutle, I., Edwards, J., Finnenkoetter, A., Franklin, C., Hanley, K., Lean,

529 H., Lock, A., Manners, J., Mittermaier, M., Morcrette, C., North, R., Petch, J., Short, C., Vosper,

530 S., Walters, D., Webster, S., ... Zerroukat, M. (2020). The first Met Office Unified Model–
531 JULES Regional Atmosphere and Land configuration, RAL1. *Geoscientific Model Development*,
532 *13*(4), 1999–2029. <https://doi.org/10.5194/gmd-13-1999-2020>

533 Chen, G., Wang, W., Yang, S., Wang, Y., Zhang, F., & Wu, K. (2023). A Neural Network-Based Scale-
534 Adaptive Cloud-Fraction Scheme for GCMs. *Journal of Advances in Modeling Earth Systems*,
535 *15*(6), e2022MS003415. <https://doi.org/10.1029/2022MS003415>

536 E3SM Project, D. (2024). *Energy Exascale Earth System Model v3.0.0* [Computer software]. [object
537 Object]. <https://doi.org/10.11578/E3SM/DC.20240301.3>

538 Gerber, R., & Wasserman, H. (2013). *High Performance Computing and Storage Requirements for*
539 *Biological and Environmental Research Target 2017* (LBNL-6256E). Lawrence Berkeley
540 National Lab. (LBNL), Berkeley, CA (United States). <https://doi.org/10.2172/1171504>

541 Gettelman, A., Gagne, D. J., Chen, C.-C., Christensen, M. W., Lebo, Z. J., Morrison, H., & Gantos, G.
542 (2021). Machine Learning the Warm Rain Process. *Journal of Advances in Modeling Earth*
543 *Systems*, *13*(2), e2020MS002268. <https://doi.org/10.1029/2020MS002268>

544 Gettelman, A., Truesdale, J. E., Bacmeister, J. T., Caldwell, P. M., Neale, R. B., Bogenschutz, P. A., &
545 Simpson, I. R. (2019). The Single Column Atmosphere Model Version 6 (SCAM6): Not a Scam
546 but a Tool for Model Evaluation and Development. *Journal of Advances in Modeling Earth*
547 *Systems*, *11*(5), 1381–1401. <https://doi.org/10.1029/2018MS001578>

548 Golaz, J.-C., Caldwell, P. M., Van Roekel, L. P., Petersen, M. R., Tang, Q., Wolfe, J. D., Abeshu, G.,
549 Anantharaj, V., Asay-Davis, X. S., Bader, D. C., Baldwin, S. A., Bisht, G., Bogenschutz, P. A.,
550 Branstetter, M., Brunke, M. A., Brus, S. R., Burrows, S. M., Cameron-Smith, P. J., Donahue, A.
551 S., ... Zhu, Q. (2019). The DOE E3SM Coupled Model Version 1: Overview and Evaluation at
552 Standard Resolution. *Journal of Advances in Modeling Earth Systems*, *11*(7), 2089–2129.
553 <https://doi.org/10.1029/2018MS001603>

554 Golaz, J.-C., Van Roekel, L. P., Zheng, X., Roberts, A. F., Wolfe, J. D., Lin, W., Bradley, A. M., Tang,
555 Q., Maltrud, M. E., Forsyth, R. M., Zhang, C., Zhou, T., Zhang, K., Zender, C. S., Wu, M.,

556 Wang, H., Turner, A. K., Singh, B., Richter, J. H., ... Bader, D. C. (2022). The DOE E3SM
557 Model Version 2: Overview of the Physical Model and Initial Model Evaluation. *Journal of*
558 *Advances in Modeling Earth Systems*, 14(12), e2022MS003156.
559 <https://doi.org/10.1029/2022MS003156>

560 Han, Y., Zhang, G. J., Huang, X., & Wang, Y. (2020). A Moist Physics Parameterization Based on Deep
561 Learning. *Journal of Advances in Modeling Earth Systems*, 12(9), e2020MS002076.
562 <https://doi.org/10.1029/2020MS002076>

563 Hartmann, D. L., Blossey, P. N., & Dygert, B. D. (2019). Convection and Climate: What Have We
564 Learned from Simple Models and Simplified Settings? *Current Climate Change Reports*, 5(3),
565 196–206. <https://doi.org/10.1007/s40641-019-00136-9>

566 Hourdin, F., Mauritsen, T., Gettelman, A., Golaz, J.-C., Balaji, V., Duan, Q., Folini, D., Ji, D., Klocke,
567 D., Qian, Y., Rauser, F., Rio, C., Tomassini, L., Watanabe, M., & Williamson, D. (2017). The Art
568 and Science of Climate Model Tuning. *Bulletin of the American Meteorological Society*, 98(3),
569 589–602. <https://doi.org/10.1175/BAMS-D-15-00135.1>

570 Huang, H., Xue, Y., Li, F., & Liu, Y. (2020). Modeling long-term fire impact on ecosystem
571 characteristics and surface energy using a process-based vegetation–fire model SSiB4/TRIFFID-
572 Fire v1.0. *Geoscientific Model Development*, 13(12), 6029–6050. [https://doi.org/10.5194/gmd-13-](https://doi.org/10.5194/gmd-13-6029-2020)
573 [6029-2020](https://doi.org/10.5194/gmd-13-6029-2020)

574 Krasnopolsky, V. M., Fox-Rabinovitz, M. S., & Belochitski, A. A. (2013). Using ensemble of neural
575 networks to learn stochastic convection parameterizations for climate and numerical weather
576 prediction models from data simulated by a cloud resolving model. *Advances in Artificial Neural*
577 *Systems*, 2013, 5–5. <https://doi.org/10.1155/2013/485913>

578 Lee, M.-I., Schubert, S. D., Suarez, M. J., Held, I. M., Lau, N.-C., Ploshay, J. J., Kumar, A., Kim, H.-K.,
579 & Schemm, J.-K. E. (2007). An Analysis of the Warm-Season Diurnal Cycle over the Continental
580 United States and Northern Mexico in General Circulation Models. *Journal of*
581 *Hydrometeorology*, 8(3), 344–366. <https://doi.org/10.1175/JHM581.1>

582 O’Gorman, P. A., & Dwyer, J. G. (2018). Using machine learning to parameterize moist convection:
583 Potential for modeling of climate, climate change, and extreme events. *Journal of Advances in*
584 *Modeling Earth Systems*, 10(10), 2548–2563. <https://doi.org/10.1029/2018ms001351>

585 Randall, D. A. (2013). Beyond deadlock. *Geophysical Research Letters*, 40(22), 5970–5976.
586 <https://doi.org/10.1002/2013GL057998>

587 Randall, D., Khairoutdinov, M., Arakawa, A., & Grabowski, W. (2003). Breaking the Cloud
588 Parameterization Deadlock. *Bulletin of the American Meteorological Society*, 84(11), 1547–1564.
589 <https://doi.org/10.1175/BAMS-84-11-1547>

590 Rasp, S., Pritchard, M. S., & Gentine, P. (2018). Deep learning to represent subgrid processes in climate
591 models. *Proceedings of the National Academy of Sciences*, 115(39), 9684–9689.
592 <https://doi.org/10.1073/pnas.1810286115>

593 Schär, C., Fuhrer, O., Arteaga, A., Ban, N., Charpiloz, C., Girolamo, S. D., Hentgen, L., Hoefler, T.,
594 Lapillonne, X., Leutwyler, D., Osterried, K., Panosetti, D., Rüdüsühli, S., Schlemmer, L.,
595 Schulthess, T. C., Sprenger, M., Ubbiali, S., & Wernli, H. (2020). Kilometer-Scale Climate
596 Models: Prospects and Challenges. *Bulletin of the American Meteorological Society*, 101(5),
597 E567–E587. <https://doi.org/10.1175/BAMS-D-18-0167.1>

598 Swann, H. (2001). Evaluation of the mass-flux approach to parametrizing deep convection. *Quarterly*
599 *Journal of the Royal Meteorological Society*, 127(574), 1239–1260.
600 <https://doi.org/10.1002/qj.49712757406>

601 Walters, D., Boutle, I., Brooks, M., Melvin, T., Stratton, R., Vosper, S., Wells, H., Williams, K., Wood,
602 N., Allen, T., Bushell, A., Copsey, D., Earnshaw, P., Edwards, J., Gross, M., Hardiman, S.,
603 Harris, C., Heming, J., Klingaman, N., ... Xavier, P. (2017). The Met Office Unified Model
604 Global Atmosphere 6.0/6.1 and JULES Global Land 6.0/6.1 configurations. *Geoscientific Model*
605 *Development*, 10(4), 1487–1520. <https://doi.org/10.5194/gmd-10-1487-2017>

606 Wang, X., Han, Y., Xue, W., Yang, G., & Zhang, G. J. (2022). Stable climate simulations using a realistic
607 general circulation model with neural network parameterizations for atmospheric moist physics

608 and radiation processes. *Geoscientific Model Development*, 15(9), 3923–3940.
609 <https://doi.org/10.5194/gmd-15-3923-2022>

610 Webster, S., Uddstrom, M., Oliver, H., & Vosper, S. (2008). A high-resolution modelling case study of a
611 severe weather event over New Zealand. *Atmospheric Science Letters*, 9(3), 119–128.
612 <https://doi.org/10.1002/asl.172>

613 Xu, K.-M., & Randall, D. A. (1996). A Semiempirical Cloudiness Parameterization for Use in Climate
614 Models. *Journal of the Atmospheric Sciences*, 53(21), 3084–3102. [https://doi.org/10.1175/1520-0469\(1996\)053<3084:ASCPFU>2.0.CO;2](https://doi.org/10.1175/1520-0469(1996)053<3084:ASCPFU>2.0.CO;2)

615

616 Zhang, T., Lin, W., Vogelmann, A. M., Zhang, M., Xie, S., Qin, Y., & Golaz, J.-C. (2021). Improving
617 Convection Trigger Functions in Deep Convective Parameterization Schemes Using Machine
618 Learning. *Journal of Advances in Modeling Earth Systems*, 13(5), e2020MS002365.
619 <https://doi.org/10.1029/2020MS002365>

620 Zhang, T., Morcrette, C., Zhang, M., Lin, W., Xie, S., Liu, Y., Weverberg, K. V., & Rodrigues, J. (2024).
621 *tzhang-ccs/ML4ESM: ML4ESM_v1* (Version v1) [Computer software]. [object Object].
622 <https://doi.org/10.5281/ZENODO.11005103>

623