



Moving beyond post-hoc XAI: Lessons learned from dynamical climate modeling

Ryan J. O'Loughlin¹, Dan Li², Travis A. O'Brien^{3,4}

¹Philosophy Department, Queens College, City University of New York, New York, 11367, USA

²Department of Philosophy, Baruch College, City University of New York, New York, 10010, USA

³Department of Earth and Atmospheric Sciences, Indiana University, Bloomington, 47405, USA

⁴Lawrence Berkeley Lab Climate and Ecosystem Sciences Division, Berkeley, CA 94720, USA

Correspondence to: Ryan J. O'Loughlin (ryan.oloughlin@qc.cuny.edu)



1 Moving beyond post-hoc XAI: Lessons learned from dynamical climate modeling

2
3

4 **Abstract.** AI models are criticized as being black boxes, potentially subjecting climate science to greater uncertainty.
5 Explainable artificial intelligence (XAI) has been proposed to probe AI models and increase trust. In this Perspective,
6 we suggest that, in addition to using XAI methods, AI researchers in climate science can learn from past successes in
7 the development of physics-based dynamical climate models. Dynamical models are complex but have gained trust
8 because their successes and failures can be attributed to specific components or sub-models, such as when model bias
9 is explained by pointing to a particular parameterization. We propose three types of understanding as a basis to
10 evaluate trust in dynamical and AI models alike: (1) instrumental understanding, which is obtained when a model has
11 passed a functional test; (2) statistical understanding, which is obtained when researchers can make sense of the
12 modelling results using statistical techniques to identify input-output relationships; and (3) Component-level
13 understanding, which refers to modelers' ability to point to specific model components or parts in the model
14 architecture as the culprit for erratic model behaviors or as the crucial reason why the model functions well. We
15 demonstrate how component-level understanding has been sought and achieved via climate model intercomparison
16 projects over the past several decades. Such component-level of understanding routinely leads to model improvements
17 and may also serve as a template for thinking about AI-driven climate science. Currently, XAI methods can help
18 explain the behaviors of AI models by focusing on the mapping between input and output, thereby increasing the
19 statistical understanding of AI models. Yet, to further increase our understanding of AI models, we will have to build
20 AI models that have interpretable components amenable to component-level understanding. We give recent examples
21 from the AI climate science literature to highlight some recent, albeit limited, successes in achieving component-level
22 understanding and thereby explaining model behaviour. The merit of such interpretable AI models is that they serve
23 as a stronger basis for trust in climate modeling and, by extension, downstream uses of climate model data.

24
25

26 **1. Introduction**

27 Machine learning (ML) is becoming increasingly utilized in climate science for tasks ranging
28 from climate model emulation (Beucler et al. 2019), to downscaling (McGinnis et al. 2021),
29 forecasting (Ham, Kim, and Luo 2019) and analyzing complex and large datasets more generally
30 (for an overview of ML in climate science, see Reichstein et al. 2019; Molina et al. 2023; de
31 Burgh-Day and Leeuwenburg 2023). Compared with physics-based methods, ML, once trained,
32 has a key advantage: computational efficiency. Along with the advantages of ML come



33 challenges such as assessing ML trustworthiness. For example, scientists often do not understand
34 why a neural net (NN) gives the output that it does because the NN is a “black box.”¹

35 To build trust in ML, the field of explainable artificial intelligence (XAI) has become
36 increasingly prominent in climate science (Bommer et al. 2023). Sometimes referred to as
37 “opening the black box,” XAI methods consist of additional models or algorithms intended to
38 shed light on why the ML model gives the output that it does. For example, (Labe and Barnes
39 2021) use an XAI method, layer-wise relevance propagation, and find that their NN heavily
40 relies on datapoints from the North Atlantic, Southern Ocean, and Southeast Asia to make its
41 predictions.

42 While XAI methods can produce useful information about ML model behaviors, these methods
43 also face problems and have been subjected to critique. As Barnes et al. (2022) note, XAI
44 methods “do not explain the actual decision-making process of the network” (p. 1). Additionally,
45 different XAI methods applied to the same ML model prediction have been shown to exhibit
46 discordance, i.e., yielding different and even incompatible “explanations” for the same ML
47 model (Mamalakis et al. 2022). Discordance in XAI is not unique to climate science. Krishna et
48 al. (2022) find that 84% of their interviewees (ML practitioners across fields who use XAI
49 methods) report experiencing discordance in their day-to-day workflow and when it comes to
50 resolving discordance, 86% of their online user study responses indicate that ML practitioners
51 either employed arbitrary heuristics (e.g., choosing a favorite method or result) or just simply did
52 not know what to do.

53 As Molina et al. (2023) note, “identifying potential failure modes of XAI, and uncertainty
54 quantification pertaining to different types of XAI methods, are both crucial to establish
55 confidence levels in XAI output and determine whether ML predictions are ‘right for the right
56 reasons’” (p. 8). Rudin (2019) argues that, instead of attempting to use XAI to explain ML
57 models post hoc, scientists ought to build interpretable models informed by domain expertise
58 from the outset. Speaking about explainability in particular, Rudin writes, “many of the [XAI]
59 methods that claim to produce *explanations* instead compute useful summary statistics of
60 predictions made by the original model. Rather than producing explanations that are faithful to
61 the original model, they show trends in how predictions are related to the features” of the model
62 input (2019, p. 208).

63 Regardless, XAI methods will likely continue to be widely applied due to ease of use and as
64 benchmark metrics for XAI methods are proposed and implemented (Hedström et al. 2023;
65 Bommer et al. 2023). In some cases, XAI methods are applied with great success, e.g.,
66 (Mamalakis et al. 2022) found that the input x gradient method fit their ground truth model with

¹ Note that computer scientists have proposed various conceptual approaches to articulate “transparency” (e.g., Lipton 2016). However, we aim to offer conceptual clarity for ML applications specifically in climate science by comparing different types of understanding in ML and in dynamical climate models.



67 a high degree of accuracy. However, we believe that more progress can be made in establishing
68 trust in ML-driven climate science.

69 In this Perspective, we recommend that climate scientists move beyond traditional post hoc XAI
70 methods and aim for *component-level* understanding of ML models. By “component” we mean a
71 functional unit of the model’s architecture, such as a layer or layers in a neural net. By
72 “understanding” we mean knowledge that could serve as a basis for an explanation about the
73 model. We distinguish between three levels of understanding:

74 **Instrumental understanding:** knowing *that* the model performed well (or not); e.g.,
75 knowing its error rate on a given test.

76 **Statistical understanding:** being able to offer a reason why we should trust a given ML
77 model by appealing to input-output mappings. These mappings can be retrieved by
78 statistical techniques.

79 **Component-level understanding:** being able to point to specific model components or
80 parts in the model architecture as the cause of erratic model behaviors or as the crucial
81 reason why the model functions well.

82 Instrumental understanding, while clearly necessary, is fairly straightforward and is a
83 prerequisite for any explanation of model behavior. It involves knowing the degree to which a
84 model fits some data (Lloyd 2010; Baumberger et al. 2017). It may also involve knowing
85 whether the model both fits some data *and* agrees with simpler models about a prediction of
86 interest or whether the model has performed well on an out-of-sample test (e.g., (Hausfather et
87 al. 2020) or according to other metrics (e.g., Gleckler et al. 2008).

88
89 However, in this perspective, we will only focus on the other two types of understanding.
90 Statistical understanding can be gained via traditional XAI methods but does not require
91 knowledge of the model’s innerworkings, i.e., its components and/or architecture (see Sect. 2
92 below). In contrast, component-level understanding *does* involve knowledge of the model’s
93 innerworkings. Therefore, component-level understanding allows scientists to offer causal
94 explanations that attribute ML model behaviors to its components. Scientists need to build and
95 analyze their models in such a way that they can understand how distinct model components
96 contribute to the model’s overall predictive successes or failures rather than merely probe model
97 data to yield input-output mappings. The latter is emblematic of traditional XAI methods.

98 Our recommendation to strive for component-level understanding is inspired by how dynamical
99 climate models have been built, tested, and improved, such as those in the coupled model
100 intercomparison projects (CMIP). In CMIP, when models agree on a particular result, scientists
101 sometimes infer that the governing equations and prescribed forcings shared by the models are
102 responsible for the models’ similar results. As Baumberger et al. (2017) put it, “robustness of
103 model results (combined with their empirical accuracy) is often seen as making it likely, or at

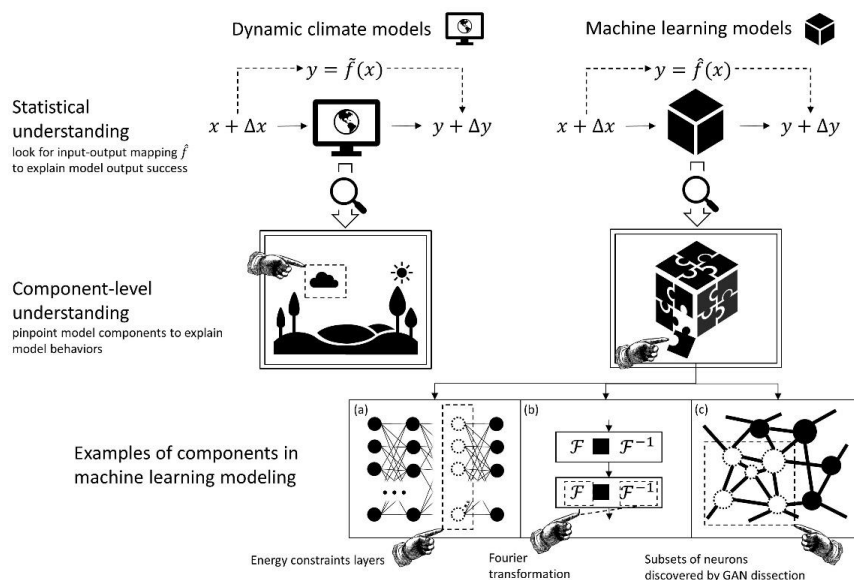


104 least increasing our confidence, that the processes that determine these results are encapsulated
105 sufficiently well in the models” (p. 11; see also Hegerl et al. 2007; Kravitz et al. 2013; Lloyd
106 2015; Schmidt and Sherwood 2015; O’Loughlin 2021). Conversely, when climate models exhibit
107 biases or errors, scientists can often point to specific parameterizations or sub-models as the
108 likely cause (e.g., Gleckler et al. 1995; Pitari et al. 2014; Gettelman et al. 2019); O’Loughlin
109 2023), although models can get the right answer for the wrong reasons (e.g., see Knutti 2008).

110 Fortunately, we see component-level understanding exemplified in ML-driven climate science to
111 some extent already (Beucler et al. 2019; Kashinath et al. 2021; Bonev et al. 2023, see Sect. 4
112 below). Indeed, the thinking behind physics-informed machine learning, which incorporates
113 known physical relations into the models from the outset (Kashinath et al. 2021; Wang et al.
114 2022; Cuomo et al. 2022), often involves component-level understanding. Thus, our proposal is
115 an endorsement of these ongoing best practices, a recognition of the relationship between the
116 evaluation of dynamical models and data-driven models, and a warning about the limits of
117 statistical understanding.

118 In addition, there is a concurrent need to establish the trustworthiness of ML models as ML-
119 driven climate science potentially becomes increasingly used to inform decision makers. While
120 decision makers themselves do not need to understand exactly how a model arrives at the answer
121 it does, they may desire an explanation of the model’s behavior that comes from a credible
122 expert. One way to establish credibility is to be able to explain ML model behavior by appealing
123 to the innerworkings of the model, which requires component-level understanding of the model.
124 In this way, component-level understanding can serve as a basis for trust in ML-driven climate
125 science.

126 The remainder of the paper is structured as follows. In Sect. 2, we give an overview of XAI in
127 climate science and explain the idea of statistical understanding and how XAI can only give us
128 statistical understanding. In Sect. 3, we detail the notion of component-level understanding and
129 demonstrate it using examples from CMIP. In Sect. 4, we show how component-level
130 understanding is achievable in ML. In Sect. 5, we conclude and make suggestions for ML-driven
131 climate science.



132

Figure 1. Scientists can obtain statistical understanding of models by seeking input-output mapping, e.g., via perturbation experiments. To acquire component-level understanding, one needs to be able to pinpoint specific components to explain models’ erratic behaviors or successes. This has been done in dynamic climate modeling, e.g., by pointing to cloud parameterization as a means to improve modeling outcomes. We offer three examples of component-level understanding in machine modeling. In panel (a), Beucler et al. (2021) design layers of neurons in their neural network to enforce energy conservation and improved model outcome. In panel (b), Kathnash et al. (2023) use spherical Fourier transformation to ensure Fourier Neural Operators perform with climate data. In panel (c), Bau et al. (2019) use a method called GAN dissection to identify which subsets of neurons control parts of images that correspond to semantics (e.g., trees or doors).

133

134 2. Post-hoc XAI in climate science and statistical understanding

135 XAI methods are intended to shed light on the behavior of complex opaque ML models. As
 136 Mamalakis et al. (2022b) put it, XAI “methods aim at a post hoc attribution of the NN prediction
 137 to specific features in the input domain (usually referred to as attribution/relevance heatmaps),
 138 thus identifying relationships between the input and the output that may be interpreted physically
 139 by the scientists” (p. 316). XAI methods are typically applied to ML models which are multi-
 140 layer, convolutional, recurrent neural networks, and/or tree ensembles.

141 The general idea behind XAI methods is to attribute the predictive success of the model’s output
 142 (i.e., the model’s prediction or decision) to subsets of its input in supervised ML. Broadly, there
 143 are two conceptual approaches to achieve this.² One approach is to figure out how the changes in
 144 input affect the output. For example, Local Interpretable Model-agnostic Explanation (LIME)
 145 first perturbs an input data point to create surrogate data near said data point. Then, after the

² Yuan et al. (2023) break down the various XAI methods into four categories. They divide those related to manipulating input-output into perturbation-based methods and surrogate-based methods (e.g., LIME). They divide the methods that rely on model parameter values into gradient-based methods (e.g., gradient) and decomposition-based method (e.g., LRP).



146 trained ML model classifies the surrogate data, LIME fits a linear regression using classified
147 surrogate data and measures how model output can be attributed to features of the surrogate data
148 manifold. In this way, LIME attributes the predictive success for the actual data point to a subset
149 of input features. Note that L stands for “local” because LIME starts with perturbing specific
150 classificatory instances rather than with global classification.

151 Another commonly used method is Shapley Additive explanation (SHAP), which is based on
152 calculating the Shapley values of each input feature. Shapley values are cooperative game
153 theoretic measures that distribute gains or costs to members of a coalition. Roughly put, Shapley
154 values are calculated by repeatedly randomly removing a member from the group to form a new
155 coalition and calculating the consequent gains and then averaging all marginal contributions to
156 all possible coalitions. In the XAI context, input features will have different Shapley values,
157 denoting their different contribution to the model’s predictive success. E.g., see (Chakraborty et
158 al. 2021; Felsche and Ludwig 2021; Cilli et al. 2022; Clare et al. 2022; Grundner et al. 2022; W.
159 Li et al. 2022; Xue et al. 2022)

160 Another approach relies on treating a trained black box model as a function to understand how
161 the input-output mapping relationship is represented by this function. For example, vanilla
162 gradient (also known as saliency) is an XAI method that relies on calculating the gradient of
163 probabilities of output being in each possible category with respect to its input and
164 backpropagates the information to its input. In this way, vanilla gradient quantifies the relative
165 importance of each element of the input vector with respect to the output, thereby attributing the
166 predictive success to subsets of input. E.g., see Balmaceda-Huarte et al. 2023; Liu et al. 2023; He
167 et al. 2024.³

168 Let’s examine how XAI methods yield statistical understanding in a detailed example. González-
169 Abad et al. (2023) use the saliency method to examine input-output mappings in three different
170 convolutional neural nets (CNNs) which were trained and used to downscale climate data. They
171 computed and produced accumulated saliency maps which account for “the overall importance
172 of the different elements” of the input data for the model’s prediction (p. 8). One of their results
173 is that, in one of the CNNs, air temperature (at 500hPa, 700 hPa, 850hPa, and 1000 hPa)
174 accumulates the highest relevance for predicting North American near-surface air temperature,
175 although different regions are apparently more relevant than others to the models’ predictions
176 (see their figure 6, p. 12). In other words, it appeared that the CNN had correctly picked up on a
177 relationship between coarse resolution temperature at certain geopotential heights on the one
178 hand, and higher resolution near-surface air temperatures on the other hand.

179 In this way, XAI methods yield information that can be helpful for making a model’s results
180 intelligible. E.g., it puts a scientist in the position to say, “this model was picking up on aspects

³ Yet another commonly used XAI method, layerwise relevance propagation (LRP), computes how each neuron contributes to other neurons’ activations, therefore highlighting the subsets of the input that dominantly contribute to the output. E.g., see (Gordon, Barnes, and Hurrell 2021; Toms, Barnes, and Hurrell 2021; Labe and Barnes 2021; 2022a; 2022b; Rader et al. 2022; Diffenbaugh and Barnes 2023).



181 A, B, and C of the input data. These aspects contributed to prediction X, a prediction that seems
182 plausible.” This exemplifies what we call “statistical understanding”, i.e., being able to offer a
183 reason why we should trust a given ML model by appealing to statistical mappings between
184 input and output. Statistical techniques are often used to obtain these mappings by relating
185 variations in input to variations in output. Post hoc XAI methods can typically yield this type of
186 understanding. Note that this is not the same as explaining the innerworkings of the model itself,
187 or what we call “component-level understanding,” because the explanation does not attribute the
188 model behaviors to ML model components, but rather is focused on input-output mapping.

189 While XAI methods can give statistical understanding of model behaviors, this type of
190 understanding has limitations. The general limitation is a familiar one, i.e., that “while XAI can
191 reveal correlations between input features and outputs, the statistics adage states: ‘correlation
192 does not imply causation’” (Molina et al. 2023, p. 8)⁴. Even if genuine causal relationships
193 between input and output can be established, we still do not know how the ML model produces a
194 certain set of output. To answer this question, ideally, we would like to know the causal role
195 played by (at least) some of the components making up the model. We would like to know about
196 at least some processes, mechanisms, constraints, or structural dependencies inside of the model,
197 rather than merely probing the ML-model-as-black-box from the outside and post hoc. While
198 XAI methods can yield information that seems plausible and physically meaningful, this
199 information may be irrelevant with respect to how the model actually arrived at a given decision
200 or prediction (Rudin 2019; Baron 2023). This, in turn, can undermine our trust in the model for
201 future applications. In contrast, with component level understanding, the causal knowledge is
202 more secure and can also inform future development and improvement of the model in question
203 and ML models in general.

204

205 **3. Understanding and Intelligibility in CMIP**

206 Dynamical models are complex but have gained trust because their successes and failures can
207 regularly be attributed to specific components or sub-models, such as when model bias is
208 explained by pointing to a particular parameterization. Indeed, the practice of diagnosing model
209 errors pre-dates the Atmospheric Model Intercomparison Project (AMIP; Gates 1992). For
210 example, differences in the representation both of radiative processes and of atmospheric
211 stratification at the poles were featured in an evaluation of why 1-D models diverged from a
212 GCM in their estimate of climate sensitivity (see Schneider 1975).

213 Later, in one of the diagnostic subprojects following AMIP, Gleckler et al. (1995) attributed
214 incorrect calculations of ocean heat transport to the models’ representations of cloud radiative
215 effects. They first found that the models’ implied ocean heat transport was partially in the wrong

⁴ To be more precise, we interpret this quote as saying that correlation does not (logically) entail causation. Correlation may be a sign that there is a causal relation in play, and correlations between events often lead us to try and relate events causally.



216 direction—northward in the Southern Hemisphere. They inferred that cloud radiative effects
217 were the culprit, explicitly noting that atmospheric GCMs at the time of their writing were
218 “known to disagree considerably in their simulations of the effects of clouds on the Earth’s
219 radiation budget (Cess et al. 1989), and hence the effects of simulated cloud-radiation
220 interactions on the implied meridional energy transports [were] immediately suspect” (Gleckler
221 et al. 1995, p. 793). They recalculated ocean heat transport using a hybrid of model data and
222 observational data. When they did this, they fixed the error—ocean heat transport turned
223 poleward. The observational data used to fix the error were of cloud radiative effects. In other
224 words, they substituted the output data linked to the problematic cloud parameterizations (a
225 *component* of the models) with observational data of cloud radiative effects. This substitution
226 resulted in a better fit with observations of and physical background knowledge of ocean heat
227 transport.

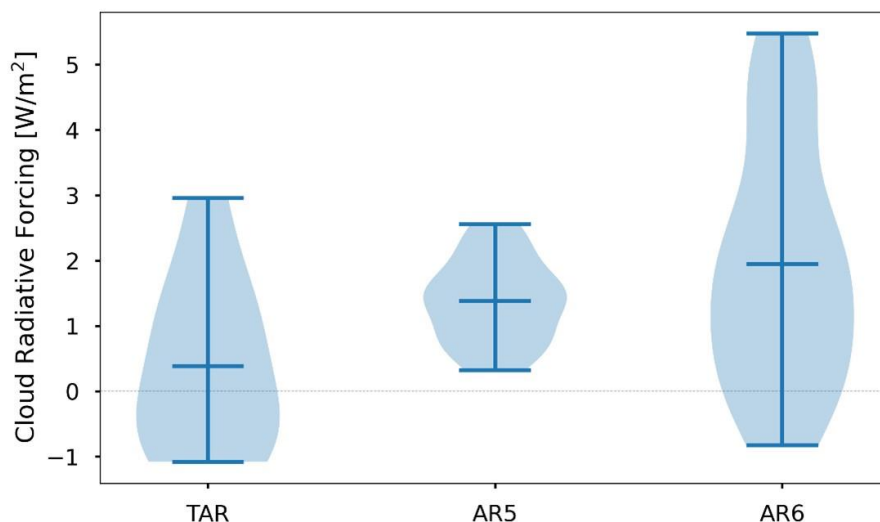
228 One may argue that substituting model components merely exemplifies statistical understanding
229 because it concerns the input and output data of the models, which, in Glecker et al.’s case, are
230 cloud-radiation and ocean heat transport. Yet, this would be misguided. Gleckler et al. isolated
231 the cloud components as the causal culprit behind why the models produced biased ocean heat
232 transport data. There is also a physically intelligible link between cloud radiative forcing and
233 ocean surface heat, so the diagnosis made scientific sense. In this way, scientists can diagnose
234 and fix climate models.

235 Many more recent cases of error diagnosis also aim to identify problematic parameterizations
236 (e.g., see (Hall and Qu 2006; O’Brien et al. 2013; Pitari et al. 2014; Bukovsky et al. 2017;
237 Gettelman et al. 2019); but see Neelin et al. 2023 for current challenges). In CMIP6 in particular,
238 there is an increased focus on process-level analysis (Eyring et al. 2019; Maloney et al. 2019). In
239 process-level analysis, scientists examine bias in the simulation of particular processes which
240 are, in turn, linked to one or more parameterizations, i.e., components within a whole GCM.
241 Moreover, CMIP-endorsed model intercomparison projects (MIPs) also center on particular
242 processes or parameterizations, such as the cloud feedbacks MIP (Webb et al. 2017) and the land
243 surface, snow and soil moisture MIP (van den Hurk et al. 2016).

244
245 The practice of updating model parameterizations during model development also demonstrates
246 an interest (and success) in achieving component-level understanding. We provide two examples
247 here: one associated with the radiative transfer parameterization in the Community Atmosphere
248 Model and another associated with the physical representation of stratocumulus clouds in
249 boundary layer parameterizations. With respect to the radiative transfer component
250 (parameterization), Collins et al. (2002) noted that, at the time their paper was written, studies
251 had “demonstrated that the longwave cooling rates and thermodynamic state simulated by GCMs
252 are sensitive to the treatment of water vapor line strengths.” Collins et al. used this knowledge—
253 along with updated information about absorption and emission of thermal radiation by water



254 vapor—to update the radiation parameterization in the Community Atmosphere Model. This
255 component-level improvement led to substantial improvements in the models’ simulated climate.



256

Figure 2. Changes in the distribution of estimated cloud radiative forcing (CRF) across three generations of IPCC Assessment Reports: 3 (TAR, published in 2001), 5 (AR5, 2014), and 6 (AR6, 2021). AR4 is omitted because data necessary to estimate CRF are not readily available. Estimates of simulated CRF were acquired by manual digitization of Figure 7.2 of Stocker et al. (2011) and by multiplying the equilibrium climate sensitivity and cloud feedback columns from Tables S1 and S2 of Zelinka et al. (2020). As the distribution of estimated cloud radiative forcing shifts upwards from TAR to AR5 to AR6, the figure shows that in AR5 and AR6, cloud feedbacks are largely positive. Indeed, AR6 states with high confidence that “future changes in clouds will, overall, cause additional warming” (Forster et al., 2021, p. 1022), yet it was not clear in TAR whether cloud feedbacks were positive. The increasing confidence in positive cloud feedbacks is partially due to improved boundary-layer parameterization, which demonstrates modelers’ component-level understanding.

257

258 Regarding stratocumulus cloud parameterization in climate models, targeted developments
259 following the Third Intergovernmental Panel on Climate Change (IPCC) Assessment Report
260 reduced uncertainty in estimates of cloud feedbacks to the extent that the 6th IPCC Assessment
261 Report now states with high confidence that “future changes in clouds will, overall, cause
262 additional warming” (p. 1022). This systematic change in cloud radiative forcing is demonstrated
263 in Figure 2. It was not clear in the Third IPCC Assessment Report (TAR) whether cloud
264 feedbacks were positive or negative, and the TAR noted in particular that the “difficulty in
265 simulation of observed boundary layer cloud properties is a clear testimony of the still
266 inadequate representation of boundary-layer processes” ((TAR 2001), p. 273). Around this time,
267 researchers developed improved boundary layer parameterizations with the goal of improving
268 the representation of low, boundary layer clouds. For instance, Grenier and Bretherton built on a
269 standard 1.5-order boundary layer turbulence parameterization in which turbulent mixing is
270 treated as a diffusive process related to the amount of turbulent kinetic energy (TKE) and in



271 which TKE is treated as a conservative, prognostic quantity. Their key additions to the 1.5-order
272 turbulence approach were (1) a more accurate numerical treatment of diffusion in the vicinity of
273 step-function-like jumps in temperature and humidity (inversions) and (2) contribution of cloud-
274 top radiative cooling to the production of TKE. These two ingredients allow the turbulence
275 parameterization to emulate the physics that drive stratocumulus clouds. Variations on the
276 parameterization of (Grenier and Bretherton (2001) and other similarly sophisticated boundary
277 layer parameterizations have been included in numerous weather and climate models, leading to
278 improvements in the simulation of stratocumulus clouds specifically and general improvements
279 in model climatology.

280 We take the above cases from CMIP to indicate that climate scientists aim for component-level
281 understanding of their models, which relates to a standard that climate models be at least
282 somewhat *intelligible*. Adopting the idea of “intelligibility” from philosopher of science Regt
283 (2017) we can say that a complex model is intelligible for scientists if they can recognize
284 qualitatively characteristic consequences of the model without performing exact calculations.
285 Intelligibility is facilitated by having models made up of components. In dynamical models,
286 these components represent real-world processes, even in cases of empirically based
287 parameterizations. More generally, knowing that a model component plays a particular role—
288 either representing the process as designed or a role later discovered during model
289 development—in a climate simulation is invaluable for reasoning about the behavior, successes,
290 and biases of the GCM as a whole.

291 The climate modeling community has long strived for component-level understanding and
292 intelligibility. This is especially evident in the work on climate model hierarchies, i.e., a group of
293 models which spans a range of complexity and comprehensiveness Jeevanjee et al. (2017).
294 Writing nearly two decades ago, Issac Held (2005) identified model hierarchies as necessary if
295 we wish to understand both the climate system and complex climate models:

296 we need a model hierarchy on which to base our understanding, describing how the dynamics
297 change as key sources of complexity are added or subtracted... (p. 1609)

298 ...the construction of such hierarchies must, I believe, be a central goal of climate theory in
299 the twenty-first century. There are no alternatives if we want to understand the climate
300 system and our comprehensive climate models. Our understanding will be embedded within
301 these hierarchies. (p. 1610)

302

303 Along similar lines, and before the advent of CMIP, Stephen Schneider (1979) wrote that

304 ...the field of climate modeling needs to “fill in the blanks” at each level in the hierarchy of
305 climate models. For only when the effect of adding one change at a time in models of
306 different complexity can be studied, will we have any real hope of understanding cause and
307 effect in the climatic system. (p. 748)



308

309 These appeals to climate model hierarchies highlight how component-level understanding is a
310 longstanding goal in climate modeling. This is not to say that component-level understanding
311 automatically translates to understanding all model behaviors. Emergent properties such as
312 equilibrium climate sensitivity may elude explanation—even when components such as cloud
313 parameterization are appealed to as causally relevant for higher ECS values (e.g., Zelinka et al.
314 2020), it must be granted that these cloud parameterizations *interact with* other components and
315 pieces of the overall GCM. So there may be a more complete explanation detailing how, as a
316 whole, the GCM simulates a higher ECS. Therefore, we do not regard our three proposed types
317 of understanding as exhaustive—perhaps a component-interaction or structural type of
318 understanding ought to be theorized and strived for as well.

319 However, the examples from earlier in this section show how the goal of component-level
320 understanding is regularly achieved, overall model complexity notwithstanding. Having achieved
321 such understanding, scientists can be more confident that their models have indeed captured
322 some truths about the target systems, and they are thereby justified to increase their confidence in
323 these complex models. In the climate modeling literature, component-level understanding
324 routinely leads to model improvements.

325 We end this section with a brief discussion distinguishing between component-level and
326 statistical understanding.

327 In general, statistical understanding can help us answer questions such as “do the input-output
328 relations of the model make sense and, if so, in what way do they make sense?” This is great for
329 finding out whether the model’s behavior is consistent with expectations across a variety of
330 cases. However, this is distinct from learning *why* the model behaves the way it does. To answer
331 this distinct question, we need to know how the model is working, which, in turn, involves
332 knowing something about the pieces making up the model. Hence, component-level
333 understanding is called for. This is exactly the type of understanding that we see aimed for, and
334 often grasped, in CMIP experiments.

335 Component-level understanding often involves a different kind of knowledge related to model
336 architecture and beyond input-output relationships. On the one hand it can demonstrate that you
337 know what role the component is playing in the model—this shows some knowledge of model-
338 building. It may also be helpful for answering a wider range of what-if-things-had-been-different
339 questions. Finally, and potentially the clearest benefit of component-level understanding, is that
340 it can tell one what needs to be fixed in cases of error. This should produce additional trust in the
341 modeling enterprise more generally.⁵

⁵ This is not to say that component-level understanding is necessarily superior to statistical understanding. E.g., knowing about a robustly detected statistical relationship could be more valuable than knowing how a single model component functions, especially since many important model behaviors arise from interactions between multiple model components.



342 4. Lessons learned: examples of component level understanding in ML

343 Component-level understanding is not the privilege solely of dynamic climate modeling. ML
344 models can be built with intelligible components as well, although their components look very
345 different from those in dynamic models. In this section, we offer three examples in which ML
346 researchers are able to acquire component-level understanding of model behaviors by
347 intentionally designing or discovering model components that are interpretable and intelligible.

348 4.1 *Attributing model success with physics-informed machine learning*

349 Our first example involves physics-informed machine learning, i.e., machine learning
350 incorporated with domain knowledge and physical principles (Kashinath et al. 2019). Model
351 success can be attributed to a specific component in a neural net, if it is known that said
352 component in the neural net is performing a physically relevant role for a given modeling task.

353 Beucler et al. (2019; 2021) augment a neural net’s architecture via layers which enforce
354 conservation laws that are important for emulating convection (see Figure 1, panel a). These laws
355 include enthalpy conservation, column-integrated water conservation, and both long- and short-
356 wave radiation conservation. The conservation laws are enforced “to machine precision”
357 (Beucler et al. 2021). Following Beucler et al. (2019) and because this neural net has a physics-
358 informed *architecture*, we will use the acronym NNA. NNA is trained on aqua-planet simulation
359 data from the Super-Parameterized Community Atmosphere Model 3.0. NNA’s results are
360 compared with those of two other neural nets: one *unconstrained* by physics (NNU) and another
361 “softly” constrained through a penalization term in the *loss* function (NNL; see Beucler et al.
362 (2019) for further discussion).

363 All three NNs are evaluated based on the mean squared errors (MSE) of their predictions and
364 based on whether their output violates physics conservation laws (P-score). While NNU has the
365 highest performance in a baseline climate—i.e., a climate well-represented by the training data—
366 NNA and NNL each outperform NNU in a 4k warmer climate (see Beucler et al. 2019, Table 1),
367 which is impressive since generalizing into warmer climate is particularly challenging for ML
368 models (Rasp et al. 2018; Li 2023). These results may indicate that NNU performed better in the
369 baseline climate for the “wrong” reasons. Indeed, NNU had a far lower P-score in both the
370 baseline and the 4k warmer climate cases.

371 Beucler et al. (2021) further show that NNA predicts the total thermodynamic tendency in the
372 enthalpy conservation equation more accurately than the other NNs—“by an amount closely
373 related to how much each NN violates enthalpy conservation” (p. 5). The particular layer in
374 NNA responsible for enthalpy conservation is obviously the explanation for this result. This case
375 therefore exemplifies component-level understanding straightforwardly.

376 It should be noted that both NNA and NNL perform well in the 4k warmer climate and, more
377 generally, “[e]nforcing constraints, whether in the architecture or the loss function, can
378 systematically reduce the error of variables that appear in the constraints” (Beucler et al. 2021, p.



379 5). This suggests that, when thinking purely about model performance, physical constraints do
380 not necessarily need to be implemented *in* the model’s architecture. However, compared with
381 NNL, Beucler et al.’s use of NNA facilitates straightforward component-level understanding.
382 The component-level understanding is straightforward because we know that, by virtue of the
383 physics knowledge built into the model’s architecture, NNA obeys conservation laws as it is
384 trained and as it is tested. We can draw an analogy with dynamical climate models. NNL is to
385 NNA as bias-corrected GCM simulations are to ones which capture relevant physical processes
386 with high-fidelity to begin with. Knowing that a model produces a physically consistent answer
387 for physical reasons is a stronger basis for trust than merely knowing that a model produces
388 physically consistent answers due to post-hoc bias correction.

389

390 *4.2 Explaining model error in a case of Fourier Neural Operators*

391 Another example involves a recent development in using machine learning to solve partial
392 differential equations: the Fourier neural operator (FNO) pioneered by Li et al. (2021). The
393 innovation of FNO is the application of Fourier transforms to enable CNN-based layers that learn
394 “solution operators” to PDEs in a scale-invariant way. Building on Li et al. (2021) demonstrated
395 that training an FNO network on output from a numerical weather prediction (NWP) model
396 produced a machine learning model that emulates NWP models with high fidelity and efficiency.
397 A key challenge noted by Pathak et al., however, was a numerical instability that limited
398 application of the FNO model to forecasts of lengths less than 10 days.

399 Analysis of the instability ultimately led the group to hypothesize that the instability was due to a
400 specific component of the FNO model: the Fourier transform itself. The problem they identified
401 was that the sine/cosine functions employed in Fourier transforms are the eigenfunctions of the
402 Laplace operator on a doubly-periodic, Euclidean geometry, whereas the desired problem (i.e.,
403 NWP) is intrinsic to an approximately spherical geometry. In essence, the Earth’s poles represent
404 a singularity that Fourier transforms on a latitude-longitude grid are not well-equipped to handle.
405 Bonev et al. (2023) adapt the FNO approach to spherical geometry by utilizing spherical
406 harmonic transforms with the Laplace-operator eigenfunctions for spherical geometries as basis
407 functions, in lieu of Fourier transforms. These eigenfunctions, the spherical harmonic functions,
408 smoothly handle the poles as a natural part of their formulation. Bonev et al. (2023) report that
409 the application of spherical harmonic transforms, rather than Fourier transforms, results in a
410 model that is numerically stable up to at least $O(100)$ days and possibly longer.

411 The application of spherical transformations stabilizes the FNO model. Bonev et al. were able to
412 fix the FNO because they could pinpoint the Fourier transformations, a component of the FNO
413 model, demonstrating scientists’ component-level understanding.⁶

⁶ Fourier transformations turn out to be useful in other contexts of ML-driven climate science because scientists can use them to understand neural networks behaviors as combinations of filters, e.g., (Subel et al. 2023).



414 *4.3 GAN dissect for future applications in ML-driven climate science*

415 The final example comes from generative adversarial networks (GANs) in computer vision. Bau
416 et al. (2018) identify particular units (i.e., sets of neurons and/or layers) in a neural net as
417 causally relevant to the generation of particular classes within images such as doors on churches.
418 They demonstrate that these units *are* actually causally relevant by showing what happens when
419 said units are ablated (essentially setting them to 0).

420 The example demonstrates component-level understanding because the units in question are
421 manipulated. Components within the architecture of the model are turned on and off and the
422 resultant effects are observed.⁷ This puts us in a position to say, for example, “these neurons are
423 responsible for generating images of trees, and we know this because turning more of these
424 neurons on yields an image with more trees (or bigger trees) and vice versa. Moreover, the other
425 aspects of the image are unchanged no matter what we do to these neurons.” Bau et al. (2018)
426 also show that visual artefacts are causally linked to particular units and can be removed using
427 this causal knowledge.

428 This case is analogous to the study from Gleckler et al. (1995) as described in Sect. 3 above.
429 Recall that the cloud radiative effects from the GCMs were “turned off” (substituted out and
430 replaced with observational data) and the calculations of ocean heat transport improved.
431 Scientists can make sense of model error because they know that a certainty deficiency in GCMs,
432 at the time, involved components of the GCMs responsible for representing clouds. In the same
433 way, Bau et al. (2018) are able to intervene on generations of images by linking units in their
434 model to particular types of image classes and examining what happens to the overall image
435 when these units are manipulated.

436 While GAN dissect isn’t currently used in climate science research, it could be used in potential
437 future applications such as in atmospheric river detection Mahesh et al. (2023). In any case, this
438 example demonstrates yet again how component-level understanding is achievable with ML.

439

440 **5. Discussion/Recommendations for practice**

441 We have argued that component-level understanding ought to be strived for in ML-driven
442 climate science. The value of component-level understanding is especially evident in the FNO
443 problem described previously (Sect. 4.2 above). Instrumental understanding allowed the group to
444 identify a performance issue (numerical ‘issues’ in the polar regions) that led to numerical
445 instability. While the group did not employ any XAI—statistical understanding—approaches, it
446 is clear that they would have been of limited value in identifying the underlying cause of the
447 numerical instability, since XAI methods only probe input-output mappings. Ultimately the
448 problem was identified and later solved by applying component-level understanding of the FNO

⁷ As a reminder to the reader, by “component” we mean a functional unit of the model’s architecture, which includes the “units” described by (Bau et al. 2018).



449 network: knowledge that a component of the network implicitly (and incorrectly) assumed a
450 Euclidean geometry for a problem on a spherical domain.

451 However, a potential objection is that component-level understanding is unnecessary because
452 XAI methods can simply be evaluated against benchmark metrics. For example, Bommer et al.
453 (2023) propose five metrics to assess XAI methods, focusing especially on the methods' output
454 data (referred to as “explanations”). These include:

455 **Robustness** of the explanation given small perturbations to input

456 **Faithfulness**, by comparing the predictions of perturbed input and those of unperturbed input
457 to determine if a feature deemed important by the XAI method does in fact change the
458 network prediction

459

460 **Randomization**, which measures how the explanation changes by perturbing the network
461 weights, similar to the robustness metric, the thinking is that “the explanation of an input x
462 should change if the model changes or if a different class is explained” (Bommer et al.
463 (2023), p.8)

464

465 **Localization**, which measures agreement between the explanation and a user-defined region
466 of interest

467 **Complexity**, a measure of how concise the highlighted features in an explanation are, and
468 assumes that “that an explanation should consist of a few strong features” to aid
469 interpretability (Bommer et al. 2023, p. 10).

470 Insofar as the metrics are deemed desirable, we agree that such an approach could help establish
471 trust in XAI. However, we view such benchmarks as complementary to, rather than a substitute
472 for, component-level understanding. This is because benchmarks yield a sort of second-order
473 statistical understanding. That is, such metrics are largely focused on aspects of input and output
474 data produced by a given XAI method. They are, in a sense, an XXAI method, an input-output
475 mapping to help make sense of another input-output mapping.

476 Therefore, our recommendation is that ML-driven climate science strive for component-level
477 understanding. This will aid in evaluating the credibility of model results, in diagnosing model
478 error, and in model development. The clearest path to component-level understanding in ML-
479 driven climate science would likely involve climate scientists helping build the ML models that
480 are used for their research and implementing physics-based and other background knowledge to
481 whatever extent feasible (Kashinath et al. 2021; Cuomo et al. 2022). Clear standards could also
482 be developed for documenting ML architecture, training procedures, and past analyses, including
483 error diagnoses (O’Loughlin 2023). Perhaps a model intercomparison project could be developed
484 to systematically evaluate ML behavior across diverse groups of researchers. Lastly, with



485 component-level understanding as a goal to strive for, scientists can better develop hybrid
486 models where both ML and dynamic modeling components are employed.

487 Back in 2005, Held wrote that climate modeling “must proceed more systematically toward the
488 creation of a hierarchy of lasting value, providing a solid framework within which our
489 understanding of the climate system, and that of future generations, is embedded” (p. 1614). We
490 think there is a parallel need in ML-driven climate science, i.e., to develop systematic standards
491 for the use and evaluation of ML models that aid in our understanding of the climate system.
492 Striving for component-level understanding of ML models is one way to help achieve this.

493

494 **Code/data availability: No data was used or generated for this research**

495 **Author contributions:** DL conceptualized the project with assistance of RO and TO; RO wrote and prepared the
496 manuscript with writing contributions from DL and TO; DL conceptualized and created the key visualization (figure
497 1); TO conceptualized and created figure 2

498 **Funding support:** This research was supported in part by (a) the Environmental Resilience Institute, funded by
499 Indiana University's Prepared for Environmental Change Grand Challenge initiative; (b) the Andrew Mellon
500 Foundation; (c) the Professional Staff Congress, PSC-CUNY Cycle 54 Research Grant

501 **Competing interests:** At least one of the (co-)authors is a member of the editorial board of Geoscientific Model
502 Development.

503

504 **References**

505 Balmaceda-Huarte, Rocío, Jorge Baño-Medina, Matias Ezequiel Olmo, and Maria Laura Bettolli. 2023.
506 “On the Use of Convolutional Neural Networks for Downscaling Daily Temperatures over
507 Southern South America in a Climate Change Scenario.” *Climate Dynamics*, August.
508 <https://doi.org/10.1007/s00382-023-06912-6>.

509 Barnes, Elizabeth A., Randal J. Barnes, Zane K. Martin, and Jamin K. Rader. 2022. “This Looks Like That
510 There: Interpretable Neural Networks for Image Tasks When Location Matters.” *Artificial
511 Intelligence for the Earth Systems* 1 (3). <https://doi.org/10.1175/AIES-D-22-0001.1>.

512 Baron, Sam. 2023. “Explainable AI and Causal Understanding: Counterfactual Approaches Considered.”
513 *Minds and Machines* 33 (2): 347–77. <https://doi.org/10.1007/s11023-023-09637-x>.

514 Bau, David, Jun-Yan Zhu, Hendrik Strobelt, Bolei Zhou, Joshua B. Tenenbaum, William T. Freeman, and
515 Antonio Torralba. 2018. “GAN Dissection: Visualizing and Understanding Generative Adversarial
516 Networks.” arXiv. <https://doi.org/10.48550/arXiv.1811.10597>.

517 Baumberger, Christoph, Reto Knutti, and Gertrude Hirsch Hadorn. 2017. “Building Confidence in Climate
518 Model Projections: An Analysis of Inferences from Fit.” *WIREs Climate Change* 8 (3): e454.
519 <https://doi.org/10.1002/wcc.454>.



- 520 Beucler, Tom, Michael Pritchard, Stephan Rasp, Jordan Ott, Pierre Baldi, and Pierre Gentine. 2021.
521 “Enforcing Analytic Constraints in Neural Networks Emulating Physical Systems.” *Physical Review*
522 *Letters* 126 (9): 098302. <https://doi.org/10.1103/PhysRevLett.126.098302>.
- 523 Beucler, Tom, Stephan Rasp, Michael Pritchard, and Pierre Gentine. 2019. “Achieving Conservation of
524 Energy in Neural Network Emulators for Climate Modeling.” arXiv.
525 <https://doi.org/10.48550/arXiv.1906.06622>.
- 526 Bommer, Philine, Marlene Kretschmer, Anna Hedström, Dilyara Bareeva, and Marina M.-C. Höhne. 2023.
527 “Finding the Right XAI Method -- A Guide for the Evaluation and Ranking of Explainable AI
528 Methods in Climate Science.” arXiv. <https://doi.org/10.48550/arXiv.2303.00652>.
- 529 Bonev, Boris, Thorsten Kurth, Christian Hundt, Jaideep Pathak, Maximilian Baust, Karthik Kashinath, and
530 Anima Anandkumar. 2023. “Spherical Fourier Neural Operators: Learning Stable Dynamics on the
531 Sphere.” arXiv. <https://doi.org/10.48550/arXiv.2306.03838>.
- 532 Bukovsky, Melissa S., Rachel R. McCrary, Anji Seth, and Linda O. Mearns. 2017. “A Mechanistically
533 Credible, Poleward Shift in Warm-Season Precipitation Projected for the U.S. Southern Great
534 Plains?” *Journal of Climate* 30 (20): 8275–98. <https://doi.org/10.1175/JCLI-D-16-0316.1>.
- 535 Burgh-Day, Catherine Odelia de, and Tennessee Leeuwenburg. 2023. “Machine Learning for Numerical
536 Weather and Climate Modelling: A Review.” *EGUsphere*, April, 1–48.
537 <https://doi.org/10.5194/egusphere-2023-350>.
- 538 Cess, R. D., G. L. Potter, J. P. Blanchet, G. J. Boer, S. J. Ghan, J. T. Kiehl, H. Le Treut, et al. 1989.
539 “Interpretation of Cloud-Climate Feedback as Produced by 14 Atmospheric General Circulation
540 Models.” *Science* 245 (4917): 513–16. <https://doi.org/10.1126/science.245.4917.513>.
- 541 Chakraborty, Debaditya, Hakan Başağaoğlu, Lilianna Gutierrez, and Ali Mirchi. 2021. “Explainable AI
542 Reveals New Hydroclimatic Insights for Ecosystem-Centric Groundwater Management.”
543 *Environmental Research Letters* 16 (11): 114024. <https://doi.org/10.1088/1748-9326/ac2fde>.
- 544 Cilli, Roberto, Mario Elia, Marina D’Este, Vincenzo Giannico, Nicola Amoroso, Angela Lombardi, Ester
545 Pantaleo, et al. 2022. “Explainable Artificial Intelligence (XAI) Detects Wildfire Occurrence in the
546 Mediterranean Countries of Southern Europe.” *Scientific Reports* 12 (1): 16349.
547 <https://doi.org/10.1038/s41598-022-20347-9>.
- 548 Clare, Mariana C. A., Maike Sonnewald, Redouane Lguensat, Julie Deshayes, and V. Balaji. 2022.
549 “Explainable Artificial Intelligence for Bayesian Neural Networks: Toward Trustworthy Predictions
550 of Ocean Dynamics.” *Journal of Advances in Modeling Earth Systems* 14 (11): e2022MS003162.
551 <https://doi.org/10.1029/2022MS003162>.
- 552 Collins, William D., Jeremy K. Hackney, and David P. Edwards. 2002. “An Updated Parameterization for
553 Infrared Emission and Absorption by Water Vapor in the National Center for Atmospheric
554 Research Community Atmosphere Model.” *Journal of Geophysical Research: Atmospheres* 107
555 (D22): ACL 17-1-ACL 17-20. <https://doi.org/10.1029/2001JD001365>.
- 556 Cuomo, Salvatore, Vincenzo Schiano Di Cola, Fabio Giampaolo, Gianluigi Rozza, Maziar Raissi, and
557 Francesco Piccialli. 2022. “Scientific Machine Learning Through Physics-Informed Neural



- 558 Networks: Where We Are and What's Next." *Journal of Scientific Computing* 92 (3): 88.
559 <https://doi.org/10.1007/s10915-022-01939-z>.
- 560 Diffenbaugh, Noah S., and Elizabeth A. Barnes. 2023. "Data-Driven Predictions of the Time Remaining
561 until Critical Global Warming Thresholds Are Reached." *Proceedings of the National Academy of
562 Sciences* 120 (6): e2207183120. <https://doi.org/10.1073/pnas.2207183120>.
- 563 Eyring, Veronika, Peter M. Cox, Gregory M. Flato, Peter J. Gleckler, Gab Abramowitz, Peter Caldwell,
564 William D. Collins, et al. 2019. "Taking Climate Model Evaluation to the next Level." *Nature
565 Climate Change* 9 (2): 102–10. <https://doi.org/10.1038/s41558-018-0355-y>.
- 566 Felsche, Elizaveta, and Ralf Ludwig. 2021. "Applying Machine Learning for Drought Prediction in a Perfect
567 Model Framework Using Data from a Large Ensemble of Climate Simulations." *Natural Hazards
568 and Earth System Sciences* 21 (12): 3679–91. <https://doi.org/10.5194/nhess-21-3679-2021>.
- 569 Gates, W. Lawrence. 1992. "AMIP: The Atmospheric Model Intercomparison Project." *Bulletin of the
570 American Meteorological Society* 73 (12): 1962–70. [https://doi.org/10.1175/1520-0477\(1992\)073<1962:ATAMIP>2.0.CO;2](https://doi.org/10.1175/1520-0477(1992)073<1962:ATAMIP>2.0.CO;2).
- 572 Gettelman, A., C. Hannay, J. T. Bacmeister, R. B. Neale, A. G. Pendergrass, G. Danabasoglu, J.-F. Lamarque,
573 et al. 2019. "High Climate Sensitivity in the Community Earth System Model Version 2 (CESM2)." *574 Geophysical Research Letters* 46 (14): 8329–37. <https://doi.org/10.1029/2019GL083978>.
- 575 Gleckler, P. J., D. A. Randall, G. Boer, R. Colman, M. Dix, V. Galin, M. Helfand, et al. 1995. "Cloud-Radiative
576 Effects on Implied Oceanic Energy Transports as Simulated by Atmospheric General Circulation
577 Models." *Geophysical Research Letters* 22 (7): 791–94. <https://doi.org/10.1029/95GL00113>.
- 578 Gleckler, P. J., K. E. Taylor, and C. Doutriaux. 2008. "Performance Metrics for Climate Models." *Journal of
579 Geophysical Research: Atmospheres* 113 (D6). <https://doi.org/10.1029/2007JD008972>.
- 580 González-Abad, Jose, Jorge Baño-Medina, and José Manuel Gutiérrez. 2023. "Using Explainability to
581 Inform Statistical Downscaling Based on Deep Learning Beyond Standard Validation Approaches."
582 arXiv. <https://doi.org/10.48550/arXiv.2302.01771>.
- 583 Gordon, Emily M., Elizabeth A. Barnes, and James W. Hurrell. 2021. "Oceanic Harbingers of Pacific
584 Decadal Oscillation Predictability in CESM2 Detected by Neural Networks." *Geophysical Research
585 Letters* 48 (21): e2021GL095392. <https://doi.org/10.1029/2021GL095392>.
- 586 Grenier, Hervé, and Christopher S. Bretherton. 2001. "A Moist PBL Parameterization for Large-Scale
587 Models and Its Application to Subtropical Cloud-Topped Marine Boundary Layers." *Monthly
588 Weather Review* 129 (3): 357–77. [https://doi.org/10.1175/1520-0493\(2001\)129<0357:AMPPFL>2.0.CO;2](https://doi.org/10.1175/1520-0493(2001)129<0357:AMPPFL>2.0.CO;2).
- 590 Grundner, Arthur, Tom Beucler, Pierre Gentine, Fernando Iglesias-Suarez, Marco A. Giorgetta, and
591 Veronika Eyring. 2022. "Deep Learning Based Cloud Cover Parameterization for ICON." *Journal of
592 Advances in Modeling Earth Systems* 14 (12): e2021MS002959.
593 <https://doi.org/10.1029/2021MS002959>.



- 594 Hall, Alex, and Xin Qu. 2006. "Using the Current Seasonal Cycle to Constrain Snow Albedo Feedback in
595 Future Climate Change." *Geophysical Research Letters* 33 (3).
596 <https://doi.org/10.1029/2005GL025127>.
- 597 Ham, Yoo-Geun, Jeong-Hwan Kim, and Jing-Jia Luo. 2019. "Deep Learning for Multi-Year ENSO Forecasts."
598 *Nature* 573 (7775): 568–72. <https://doi.org/10.1038/s41586-019-1559-7>.
- 599 Hausfather, Zeke, Henri F. Drake, Tristan Abbott, and Gavin A. Schmidt. 2020. "Evaluating the
600 Performance of Past Climate Model Projections." *Geophysical Research Letters* 47 (1):
601 e2019GL085378. <https://doi.org/10.1029/2019GL085378>.
- 602 He, Renfei, Limao Zhang, and Alvin Wei Ze Chew. 2024. "Data-Driven Multi-Step Prediction and Analysis
603 of Monthly Rainfall Using Explainable Deep Learning." *Expert Systems with Applications* 235
604 (January): 121160. <https://doi.org/10.1016/j.eswa.2023.121160>.
- 605 Hedström, Anna, Leander Weber, Daniel Krakowczyk, Dilyara Bareeva, Franz Motzkus, Wojciech Samek,
606 Sebastian Lapuschkin, and Marina M.-C. Höhne. 2023. "Quantus: An Explainable Ai Toolkit for
607 Responsible Evaluation of Neural Network Explanations and Beyond." *Journal of Machine
608 Learning Research* 24 (34): 1–11.
- 609 Hegerl, Gabriele C, Francis W Zwiers, Pascale Braconnot, Nathan P Gillett, Yong Luo, Jose A Marengo
610 Orsini, Neville Nicholls, et al. n.d. "Understanding and Attributing Climate Change." In *Climate
611 Change 2007: The Physical Science Basis. Contribution of Working Group I to the Fourth
612 Assessment Report of the Intergovernmental Panel on Climate Change*, edited by S. Solomon, D.
613 Qin, M. Manning, Z. Chen, M. Marquis, K.B. Avery, M. Tignor, and H.L. Miller, 84. Cambridge
614 University Press (U.K.; New York).
- 615 Held, Isaac M. 2005. "The Gap between Simulation and Understanding in Climate Modeling." *Bulletin of
616 the American Meteorological Society* 86 (11): 1609–14. <https://doi.org/10.1175/BAMS-86-11-1609>.
- 618 Hurk, Bart van den, Hyungjun Kim, Gerhard Krinner, Sonia I. Seneviratne, Chris Derksen, Taikan Oki,
619 Hervé Douville, et al. 2016. "LS3MIP (v1.0) Contribution to CMIP6: The Land Surface, Snow and
620 Soil Moisture Model Intercomparison Project – Aims, Setup and Expected Outcome."
621 *Geoscientific Model Development* 9 (8): 2809–32. <https://doi.org/10.5194/gmd-9-2809-2016>.
- 622 Jeevanjee, Nadir, Pedram Hassanzadeh, Spencer Hill, and Aditi Sheshadri. 2017. "A Perspective on
623 Climate Model Hierarchies." *Journal of Advances in Modeling Earth Systems* 9 (4): 1760–71.
624 <https://doi.org/10.1002/2017MS001038>.
- 625 Kashinath, K., M. Mustafa, A. Albert, J-L. Wu, C. Jiang, S. Esmailzadeh, K. Azizzadenesheli, et al. 2021.
626 "Physics-Informed Machine Learning: Case Studies for Weather and Climate Modelling."
627 *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering
628 Sciences* 379 (2194): 20200093. <https://doi.org/10.1098/rsta.2020.0093>.
- 629 Knutti, Reto. 2008. "Why Are Climate Models Reproducing the Observed Global Surface Warming so
630 Well?" *Geophysical Research Letters* 35 (18). <https://doi.org/10.1029/2008GL034932>.



- 631 Kravitz, Ben, Alan Robock, Piers M. Forster, James M. Haywood, Mark G. Lawrence, and Hauke Schmidt.
632 2013. "An Overview of the Geoengineering Model Intercomparison Project (GeoMIP): GEOMIP
633 INTRODUCTION." *Journal of Geophysical Research: Atmospheres* 118 (23): 13,103-13,107.
634 <https://doi.org/10.1002/2013JD020569>.
- 635 Krishna, Satyapriya, Tessa Han, Alex Gu, Javin Pombra, Shahin Jabbari, Steven Wu, and Himabindu
636 Lakkaraju. 2022. "The Disagreement Problem in Explainable Machine Learning: A Practitioner's
637 Perspective." arXiv. <https://doi.org/10.48550/arXiv.2202.01602>.
- 638 Labe, Zachary M., and Elizabeth A. Barnes. 2021. "Detecting Climate Signals Using Explainable AI With
639 Single-Forcing Large Ensembles." *Journal of Advances in Modeling Earth Systems* 13 (6):
640 e2021MS002464. <https://doi.org/10.1029/2021MS002464>.
- 641 ———. 2022a. "Comparison of Climate Model Large Ensembles With Observations in the Arctic Using
642 Simple Neural Networks." *Earth and Space Science* 9 (7): e2022EA002348.
643 <https://doi.org/10.1029/2022EA002348>.
- 644 ———. 2022b. "Predicting Slowdowns in Decadal Climate Warming Trends With Explainable Neural
645 Networks." *Geophysical Research Letters* 49 (9): e2022GL098173.
646 <https://doi.org/10.1029/2022GL098173>.
- 647 Li, Dan. 2023. "Machines Learn Better with Better Data Ontology: Lessons from Philosophy of Induction
648 and Machine Learning Practice." *Minds and Machines*, June. [https://doi.org/10.1007/s11023-
649 023-09639-9](https://doi.org/10.1007/s11023-023-09639-9).
- 650 Li, Wantong, Mirco Migliavacca, Matthias Forkel, Jasper M. C. Denissen, Markus Reichstein, Hui Yang,
651 Gregory Duveiller, Ulrich Weber, and Rene Orth. 2022. "Widespread Increasing Vegetation
652 Sensitivity to Soil Moisture." *Nature Communications* 13 (1): 3959.
653 <https://doi.org/10.1038/s41467-022-31667-9>.
- 654 Li, Zongyi, Nikola Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik Bhattacharya, Andrew Stuart,
655 and Anima Anandkumar. 2021. "Fourier Neural Operator for Parametric Partial Differential
656 Equations." arXiv. <https://doi.org/10.48550/arXiv.2010.08895>.
- 657 Lipton, Zachary C. 2016. "The Mythos of Model Interpretability (2016)." *arXiv Preprint arXiv:1606.03490*.
- 658 Liu, Yumin, Kate Duffy, Jennifer G. Dy, and Auroop R. Ganguly. 2023. "Explainable Deep Learning for
659 Insights in El Niño and River Flows." *Nature Communications* 14 (1): 339.
660 <https://doi.org/10.1038/s41467-023-35968-5>.
- 661 Lloyd, Elisabeth A. 2010. "Confirmation and Robustness of Climate Models." *Philosophy of Science* 77 (5):
662 971–84. <https://doi.org/10.1086/657427>.
- 663 ———. 2015. "Model Robustness as a Confirmatory Virtue: The Case of Climate Science." *Studies in
664 History and Philosophy of Science Part A* 49 (February): 58–68.
665 <https://doi.org/10.1016/j.shpsa.2014.12.002>.
- 666 Mahesh, Ankur, Travis O'Brien, Burlen Loring, Abdelrahman Elbashandy, William Boos, and William
667 Collins. 2023. "Identifying Atmospheric Rivers and Their Poleward Latent Heat Transport with



- 668 Generalizable Neural Networks: ARCNNv1." *EGUsphere*, June, 1–36.
669 <https://doi.org/10.5194/egusphere-2023-763>.
- 670 Maloney, Eric D., Andrew Gettelman, Yi Ming, J. David Neelin, Daniel Barrie, Annarita Mariotti, C.-C.
671 Chen, et al. 2019. "Process-Oriented Evaluation of Climate and Weather Forecasting Models."
672 *Bulletin of the American Meteorological Society* 100 (9): 1665–86.
673 <https://doi.org/10.1175/BAMS-D-18-0042.1>.
- 674 Mamalakis, Antonios, Imme Ebert-Uphoff, and Elizabeth A. Barnes. 2022a. "Neural Network Attribution
675 Methods for Problems in Geoscience: A Novel Synthetic Benchmark Dataset." *Environmental
676 Data Science* 1.
- 677 Mamalakis, Antonios, Imme Ebert-Uphoff, and Elizabeth A. Barnes. 2022b. "Explainable Artificial
678 Intelligence in Meteorology and Climate Science: Model Fine-Tuning, Calibrating Trust
679 and Learning New Science." In *xxAI - Beyond Explainable AI: International Workshop, Held in
680 Conjunction with ICML 2020, July 18, 2020, Vienna, Austria, Revised and Extended Papers*, edited
681 by Andreas Holzinger, Randy Goebel, Ruth Fong, Taesup Moon, Klaus-Robert Müller, and
682 Wojciech Samek, 315–39. Lecture Notes in Computer Science. Cham: Springer International
683 Publishing. https://doi.org/10.1007/978-3-031-04083-2_16.
- 684 McGinnis, Seth, Daniel Korytina, Melissa Bukovsky, Rachel McCrary, and Linda Mearns. 2021. "Credibility
685 Evaluation of a Convolutional Neural Net for Downscaling GCM Output over the Southern Great
686 Plains" 2021 (December): GC42A-03.
- 687 Molina, Maria J., Travis A. O'Brien, Gemma Anderson, Moetasim Ashfaq, Katrina E. Bennett, William D.
688 Collins, Katherine Dagon, Juan M. Restrepo, and Paul A. Ullrich. 2023. "A Review of Recent and
689 Emerging Machine Learning Applications for Climate Variability and Weather Phenomena."
690 *Artificial Intelligence for the Earth Systems* 1 (aop): 1–46. <https://doi.org/10.1175/AIES-D-22-0086.1>.
- 692 Neelin, J. David, John P. Krasting, Aparna Radhakrishnan, Jessica Liptak, Thomas Jackson, Yi Ming,
693 Wenhao Dong, et al. 2023. "Process-Oriented Diagnostics: Principles, Practice, Community
694 Development, and Common Standards." *Bulletin of the American Meteorological Society* 104 (8):
695 E1452–68. <https://doi.org/10.1175/BAMS-D-21-0268.1>.
- 696 O'Brien, Travis A., Fuyu Li, William D. Collins, Sara A. Rauscher, Todd D. Ringler, Mark Taylor, Samson M.
697 Hagos, and L. Ruby Leung. 2013. "Observed Scaling in Clouds and Precipitation and Scale
698 Incognizance in Regional to Global Atmospheric Models." *Journal of Climate* 26 (23): 9313–33.
699 <https://doi.org/10.1175/JCLI-D-13-00005.1>.
- 700 O'Loughlin, Ryan. 2021. "Robustness Reasoning in Climate Model Comparisons." *Studies in History and
701 Philosophy of Science Part A* 85 (February): 34–43. <https://doi.org/10.1016/j.shpsa.2020.12.005>.
- 702 ———. 2023. "Diagnosing Errors in Climate Model Intercomparisons." *European Journal for Philosophy
703 of Science* 13 (2): 20. <https://doi.org/10.1007/s13194-023-00522-z>.
- 704 Pitari, Giovanni, Valentina Aquila, Ben Kravitz, Alan Robock, Shingo Watanabe, Irene Cionni, Natalia De
705 Luca, Glauco Di Genova, Eva Mancini, and Simone Tilmes. 2014. "Stratospheric Ozone Response
706 to Sulfate Geoengineering: Results from the Geoengineering Model Intercomparison Project



- 707 (GeoMIP).” *Journal of Geophysical Research: Atmospheres* 119 (5): 2629–53.
708 <https://doi.org/10.1002/2013JD020566>.
- 709 Rader, Jamin K., Elizabeth A. Barnes, Imme Ebert-Uphoff, and Chuck Anderson. 2022. “Detection of
710 Forced Change Within Combined Climate Fields Using Explainable Neural Networks.” *Journal of*
711 *Advances in Modeling Earth Systems* 14 (7): e2021MS002941.
712 <https://doi.org/10.1029/2021MS002941>.
- 713 Rasp, Stephan, Michael S. Pritchard, and Pierre Gentine. 2018. “Deep Learning to Represent Subgrid
714 Processes in Climate Models.” *Proceedings of the National Academy of Sciences* 115 (39): 9684–
715 89. <https://doi.org/10.1073/pnas.1810286115>.
- 716 Regt, Henk W. de. 2017. *Understanding Scientific Understanding*. Oxford University Press.
- 717 Reichstein, Markus, Gustau Camps-Valls, Bjorn Stevens, Martin Jung, Joachim Denzler, Nuno Carvalhais,
718 and Prabhat. 2019. “Deep Learning and Process Understanding for Data-Driven Earth System
719 Science.” *Nature* 566 (7743): 195–204. <https://doi.org/10.1038/s41586-019-0912-1>.
- 720 Rudin, Cynthia. 2019. “Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and
721 Use Interpretable Models Instead.” *Nature Machine Intelligence* 1 (5): 206–15.
722 <https://doi.org/10.1038/s42256-019-0048-x>.
- 723 Schmidt, Gavin A., and Steven Sherwood. 2015. “A Practical Philosophy of Complex Climate Modelling.”
724 *European Journal for Philosophy of Science* 5 (2): 149–69. [https://doi.org/10.1007/s13194-014-](https://doi.org/10.1007/s13194-014-0102-9)
725 0102-9.
- 726 Schneider, S. H. 1979. “Verification of Parameterizations in Climate Modeling.” In *Report of the Study*
727 *Conference on Climate Models: Performance, Intercomparison and Sensitivity Studies*, edited by
728 W. Lawrence Gates, 728–51. World Meteorological Organization, Global Atmospheric Research
729 Program, GARP Publications Series no. 22, 2 vols.
- 730 Schneider, Stephen H. 1975. “On the Carbon Dioxide–Climate Confusion.” *Journal of Atmospheric*
731 *Sciences* 32 (11): 2060–66. [https://doi.org/10.1175/1520-](https://doi.org/10.1175/1520-0469(1975)032<2060:OTCDC>2.0.CO;2)
732 0469(1975)032<2060:OTCDC>2.0.CO;2.
- 733 Subel, Adam, Yifei Guan, Ashesh Chattopadhyay, and Pedram Hassanzadeh. 2023. “Explaining the Physics
734 of Transfer Learning in Data-Driven Turbulence Modeling.” *PNAS Nexus* 2 (3): pgad015.
735 <https://doi.org/10.1093/pnasnexus/pgad015>.
- 736 TAR. 2001. “7.2.2.3 Boundary-Layer Mixing and Cloudiness from the IPCC’s Third Assessment Report
737 (TAR) Working Group 1.” 2001. <https://archive.ipcc.ch/ipccreports/tar/wg1/273.htm>.
- 738 Toms, Benjamin A., Elizabeth A. Barnes, and James W. Hurrell. 2021. “Assessing Decadal Predictability in
739 an Earth-System Model Using Explainable Neural Networks.” *Geophysical Research Letters* 48
740 (12): e2021GL093842. <https://doi.org/10.1029/2021GL093842>.
- 741 Wang, Sifan, Shyam Sankaran, and Paris Perdikaris. 2022. “Respecting Causality Is All You Need for
742 Training Physics-Informed Neural Networks.” arXiv. <https://doi.org/10.48550/arXiv.2203.07404>.



- 743 Webb, Mark J., Timothy Andrews, Alejandro Bodas-Salcedo, Sandrine Bony, Christopher S. Bretherton,
744 Robin Chadwick, H el ene Chepfer, et al. 2017. "The Cloud Feedback Model Intercomparison
745 Project (CFMIP) Contribution to CMIP6." *Geoscientific Model Development* 10 (1): 359–84.
746 <https://doi.org/10.5194/gmd-10-359-2017>.
- 747 Xue, Pengfei, Aditya Wagh, Gangfeng Ma, Yilin Wang, Yongchao Yang, Tao Liu, and Chenfu Huang. 2022.
748 "Integrating Deep Learning and Hydrodynamic Modeling to Improve the Great Lakes Forecast."
749 *Remote Sensing* 14 (11): 2640. <https://doi.org/10.3390/rs14112640>.
- 750 Yuan, Hao, Haiyang Yu, Shurui Gui, and Shuiwang Ji. 2023. "Explainability in Graph Neural Networks: A
751 Taxonomic Survey." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45 (5): 5782–
752 99. <https://doi.org/10.1109/TPAMI.2022.3204236>.
- 753 Zelinka, Mark D., Timothy A. Myers, Daniel T. McCoy, Stephen Po-Chedley, Peter M. Caldwell, Paulo
754 Ceppi, Stephen A. Klein, and Karl E. Taylor. 2020. "Causes of Higher Climate Sensitivity in CMIP6
755 Models." *Geophysical Research Letters* 47 (1): e2019GL085782.
756 <https://doi.org/10.1029/2019GL085782>.
- 757
- 758