Geoscientific
Model Development

# Moving beyond post hoc explainable artificial intelligence: a perspective paper on lessons learned from dynamical climate modeling

**Ryan J. O'Loughlin[1], Dan Li[2], Richard Neale[3], and Travis A. O'Brien[4,5]**

[1]Philosophy Department, Queens College, City University of New York, New York, NY 11367, USA
[2]Department of Philosophy, Baruch College, City University of New York, New York, NY 10010, USA
[3]National Center for Atmospheric Research, Boulder, CO 80305, USA
[4]Department of Earth and Atmospheric Sciences, Indiana University, Bloomington, IN 47405, USA
[5]Lawrence Berkeley Lab Climate and Ecosystem Sciences Division, Berkeley, CA 94720, USA

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

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

## 1 Introduction

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

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

To build trust in ML, the field of explainable artificial intelligence (XAI) has become increasingly prominent in climate science (Bommer et al., 2023). Sometimes referred to as "opening the black box," XAI methods consist of additional models or algorithms intended to shed light on why the ML model gives the output that it does. For example, Labe and Barnes (2021) use an XAI method, layer-wise relevance propagation[2], and find that their NN heavily relies on data points from the North Atlantic, Southern Ocean, and Southeast Asia to make its predictions.

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

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

trends in how predictions are related to the features" of the model input Rudin (2019, p. 208).

Regardless, XAI methods will likely continue to be widely applied due to ease of use and as benchmark metrics for XAI methods are proposed and implemented (Hedström et al., 2023; Bommer et al., 2023). In some cases, XAI methods are applied with great success; e.g., Mamalakis et al. (2022b) found that the input × gradient method fit their ground-truth model with a high degree of accuracy. However, we believe that more progress can be made in establishing trust in ML-driven climate science, especially as an increasing number of researchers start incorporating ML into climate research (see Fig. 1).

In this review and perspective paper, we target readers with expertise in traditional approaches for climate science (e.g., development, evaluation, and application of traditional Earth system models) who are starting to utilize ML in their research and who may see XAI as a tempting way to gain insight into model behavior and to build confidence. From this perspective, we draw from some ideas in the philosophy of science to recommend that such researchers leverage the expanding array of freely available ML resources to move beyond post hoc XAI methods and aim for component-level understanding of ML models. By "component" we mean a functional unit of the model's architecture, such as a layer or layers in a neural net. By "understanding" we mean knowledge that could serve as a basis for an explanation about the model. We distinguish between three levels of understanding.

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

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

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

These levels concern the degree to which complex models are intelligible or graspable to scientists (De Regt and Dieks, 2005; Knüsel and Baumberger, 2020; De Regt, 2017). Therefore, our proposal has a narrower but deeper focus than recent philosophy of science accounts of understanding climate phenomena with or by using ML and dynamical climate models (Jebeile et al., 2021; Knüsel and Baumberger, 2020). We are concerned with understanding, diagnosing, and improving model behavior to inform model development.

Instrumental understanding, while clearly necessary, is fairly straightforward and is a prerequisite for any explanation of model behavior. It involves knowing the degree to

---

[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 of ML and dynamical climate models.

[2]Layer-wise relevance propagation (LRP) is a method for interpreting neural networks by calculating how each neuron contributes to the network's output. It propagates relevance scores backward from the output to the input layer, identifying which parts of the input most significantly influence the output prediction. See, e.g., Gordon et al. (2021), Toms et al. (2021), Labe and Barnes (2021, 2022a, b), Rader et al. (2022), and Diffenbaugh and Barnes (2023).
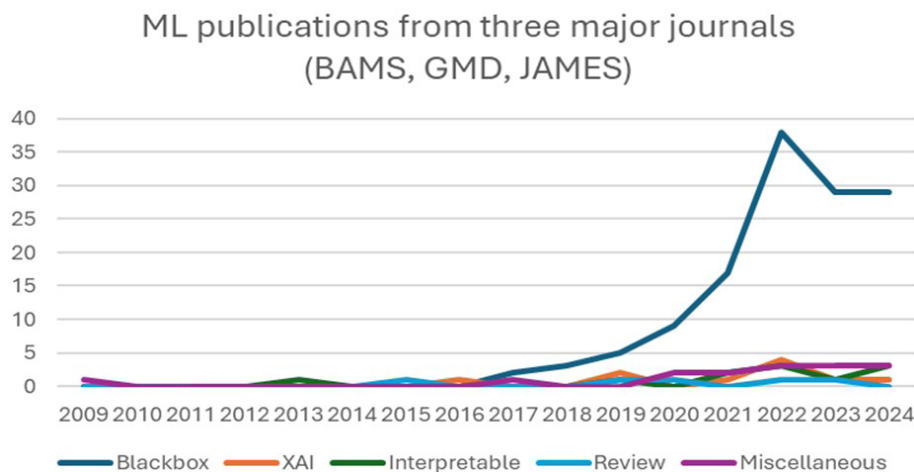
ML publications from three major journals
(BAMS, GMD, JAMES)



**Figure 1.** Trends of publications related to AI, ML, or XAI from three major journals (BAMS, GMD, JAMES) based on 178 references acquired from the Web of Science on 17 October 2023. There are two notable results. Result 1: AI and ML publications are predominantly "black box applications" (132 out of 178 records). Both XAI and interpretable AI emerged in 2020 and are in their infancy by comparison. Result 2: black box AI applications are on the rise. (The dip in 2023 and 2024 can be explained by data collection methods. Note that this is not intended as a systematic survey of AI in climate science. One key shortcoming is that the data excluded journals such as *Artificial Intelligence for Earth Systems* (AIES) because the Web of Science did not include AIES at the time of data collection.) See the Supplement for a full data description.

which a model fits some data (Lloyd, 2010; Baumberger et al., 2017). It may also involve knowing whether the model both fits some data and agrees with simpler models about a prediction of interest or whether the model has performed well on an out-of-sample test (e.g., Hausfather et al., 2020) or according to other metrics (e.g., Gleckler et al., 2008).

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

Our recommendation to strive for component-level understanding is inspired by how dynamical climate models have been built, tested, and improved, such as those in the Coupled Model Intercomparison Project (CMIP). Therefore, a novel contribution of this paper is the linking of existing climate model development practices to practices that could be employed in ML model development.

In CMIP, when models agree on a particular result, scientists sometimes infer that the governing equations and prescribed forcings shared by the models are responsible for the

models' similar results. As Baumberger et al. (2017) put it, "robustness of model results (combined with their empirical accuracy) is often seen as making it likely, or at least increasing our confidence, that the processes that determine these results are encapsulated sufficiently well in the models" (p. 11; see also Hegerl et al., 2007; Kravitz et al., 2013; Lloyd, 2015; Schmidt and Sherwood, 2015; O'Loughlin, 2021). Conversely, when climate models exhibit biases or errors, scientists can often point to specific parameterizations or sub-models as the likely cause (e.g., Gleckler et al., 1995; Pitari et al., 2014; Gettelman et al., 2019; Zelinka et al., 2020; O'Loughlin, 2023), although models can get the right answer for the wrong reasons (see, e.g., Knutti, 2008).

To be clear, there are limits to how much component-level understanding can be achieved in CMIP. Dynamical climate models exhibit fuzzy (rather than sharp) modularity, meaning that the behavior of a fully coupled model is "the complex result of the interaction of the modules – not the interaction of the results of the modules" (Lenhard and Winsberg, 2010, p. 256). Climate scientists are familiar with a related problem: the difficulty in explaining how climate models generate (or not) emergent phenomena like the Madden–Julian oscillation (Lin et al., 2024). Despite these difficulties, philosophers and other scholars of climate science have documented successes in attributing model behavior to individual model components in the climate science literature (Carrier and Lenhard, 2019; Frigg et al., 2015; Gettelman et al., 2019; Hall and Qu, 2006; Hourdin et al., 2013; Mayernik, 2021; Notz et al., 2013; O'Loughlin, 2023; Oreopoulos et al., 2012; Pincus et al., 2016; Touzé-Peiffer et al., 2020). These successes do not imply anything like a "full" or "complete"

understanding of all model behavior; rather, the component-level understanding of climate model behavior comes in degrees (Jebeile et al., 2021).

Fortunately, we see component-level understanding exemplified in ML-driven climate science to some extent already (Beucler et al., 2019; Kashinath et al., 2021; Bonev et al., 2023; see Sect. 4 below). Indeed, the thinking behind physics-informed machine learning, which incorporates known physical relations into the models from the outset (Kashinath et al., 2021; Wang et al., 2022; Cuomo et al., 2022), often involves component-level understanding. Thus, our proposal is an endorsement of these ongoing best practices, a recognition of the relationship between the evaluation of dynamical models and data-driven models, and a warning about the limits of statistical understanding. In addition, there is a concurrent need to establish the trustworthiness of ML models as ML-driven climate science potentially becomes increasingly used to inform decision-makers (NSF AI Institute for Research on Trustworthy AI in Weather, Climate, and Coastal Oceanography (AI2ES), 2024). While decision-makers themselves do not need to understand exactly how a model arrives at the answer it does, they may desire an explanation of the model's behavior that comes from a credible expert. One way to establish credibility is to be able to explain ML model behavior by appealing to the inner workings of the model, which requires component-level understanding of the model. In this way, component-level understanding can serve as a basis for trust in ML-driven climate science.

The remainder of the paper is structured as follows. In Sect. 2, we give an overview of XAI in climate science and explain the idea of statistical understanding and how XAI can only give us statistical understanding. In Sect. 3, we detail the notion of component-level understanding and demonstrate it using examples from CMIP. In Sect. 4, we show how component-level understanding is achievable in ML. In Sect. 5, we conclude and make suggestions for ML-driven climate science, including describing some resources that interested readers might utilize to build the expertise in ML model design necessary to probe, build, and adapt models in a way that is amenable to component-level understanding.

## 2  Post hoc XAI in climate science and statistical understanding

XAI methods are intended to shed light on the behavior of complex opaque ML models. As Mamalakis et al. (2022a) put it, XAI "methods aim at a post hoc attribution of the NN prediction to specific features in the input domain (usually referred to as attribution/relevance heatmaps[3]), thus identify-

ing relationships between the input and the output that may be interpreted physically by the scientists" (p. 316). XAI methods are typically applied to ML models which are multi-layer, convolutional, recurrent neural networks, and/or ensembles of decision trees (a common example of the latter is random forests).

The general idea behind XAI methods is to attribute the predictive success of the model's output (i.e., the model's prediction or decision) to subsets of its input in supervised ML. Broadly, there are two conceptual approaches to achieve this.[4] One approach is perturbing the input data to figure out how the changes in input affect the output. The other approach studies the functional representation between input and output.

For the approach of perturbing the input, Local Interpretable Model-agonistic Explanation (LIME) is a method that first perturbs an input data point to create surrogate data near said data point. Then, after the trained ML model classifies the surrogate data, LIME fits a linear regression using classified surrogate data and measures how model output can be attributed to features of the surrogate data manifold. In this way, LIME attributes the predictive success for the actual data point to a subset of input features. Note that L stands for "local" because LIME, as a method, perturbs classification instances. For example, pixels, or clusters of pixels, of one image may be perturbed to create a surrogate instance. Then this surrogate instance is classified by the ML model in question to see how the output changes. LIME does not deal with all data points all at once.

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

The other approach relies on treating a trained black box model as a function to understand how the input–output mapping relationship is represented by this function. For ex-

---

[3]In this context, *attribution/relevance heatmaps* refer to arrays (typically rendered as images) indicating the parts of an input dataset that most contribute to the ML model output. In the context of image classification, for a hypothetical example, an attribution heatmap of an image with a human hand might highlight pixels as-

sociated with fingers as being particularly important for classifying the image as a hand.

[4]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 methods (e.g., layer-wise relevance propagation).
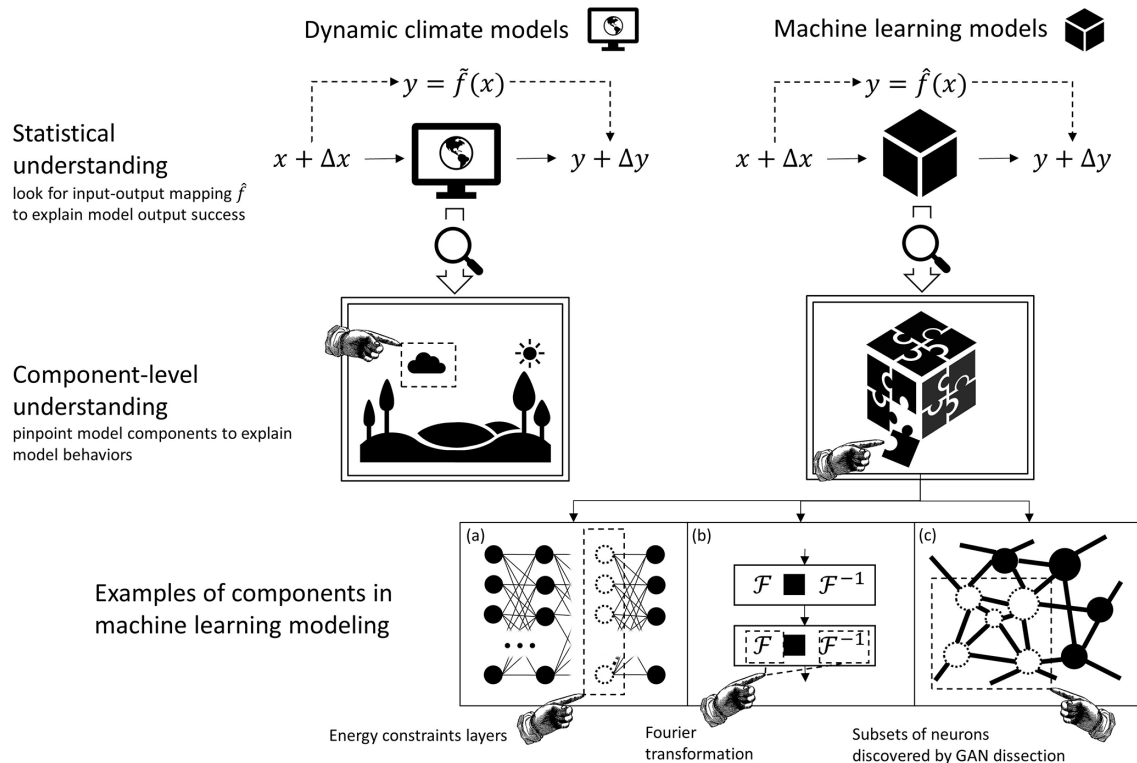
**Figure 2.** Scientists can obtain statistical understanding of models by seeking input–output mapping, e.g., via perturbation experiments. To acquire component-level understanding, one needs to be able to pinpoint specific components to explain models' erratic behaviors or successes. This has been done in dynamic climate modeling, e.g., by analyzing cloud parameterizations as a means to improve modeling outcomes. We offer three examples of component-level understanding in machine learning. In panel **(a)**, Beucler et al. (2021) design layers of neurons in their neural network to enforce energy conservation and improve model outcome. In panel **(b)**, Bonev et al. (2023) use spherical Fourier transformation to ensure that Fourier neural operators perform with climate data. In panel **(c)**, Bau et al. (2018) 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). See Sect. 4 below for a discussion of these three examples.

ample, vanilla gradient (also known as saliency) is an XAI method that relies on calculating the gradient of probabilities of output being in each possible category with respect to its input and back-propagating the information to its input. In this way, vanilla gradient quantifies the relative importance of each element of the input vector with respect to the output, thereby attributing the predictive success to subsets of input; see, e.g., Balmaceda-Huarte et al. (2023), Liu et al. (2023), and He et al. (2024).

Let us examine how XAI methods yield statistical understanding in a detailed example. González-Abad et al. (2023) use the saliency method to examine input–output mappings in three different convolutional neural nets (CNNs) which were trained and used to downscale climate data. They computed and produced accumulated saliency maps which account for "the overall importance of the different elements" of the input data for the model's prediction (p. 8). One of their results is that, in one of the CNNs, air temperature (at 500, 700, 850, and 1000 hPa) accumulates the highest relevance for predicting North American near-surface air temperature,

although different regions are apparently more relevant than others to the models' predictions (see their Fig. 6, p. 12). In other words, it appeared that the CNN had correctly picked up on a relationship between coarse-resolution temperature at certain geopotential heights on the one hand and higher-resolution near-surface air temperatures on the other hand.

In this way, XAI methods yield information that can be helpful in making a model's results intelligible. For example, it puts a scientist in the position to say, "this model was picking up on aspects A, B, and C of the input data. These aspects contributed to prediction X, a prediction that seems plausible." This exemplifies what we call "statistical understanding", i.e., being able to offer a reason why we should trust a given ML model by appealing to statistical mappings between input and output. Statistical techniques are often used to obtain these mappings by relating variations in input to variations in output. Post hoc XAI methods can typically yield this type of understanding. Note that this is not the same as explaining the inner workings of the model itself, or what we call "component-level understanding", because the expla-

nation does not attribute the model behaviors to ML model components but is rather focused on input–output mapping.

While XAI methods can give statistical understanding of model behaviors, this type of understanding has limitations. The general limitation is a familiar one, i.e., that "while XAI can reveal correlations between input features and outputs, the statistics adage states: 'correlation does not imply causation'" (Molina et al., 2023, p. 8)[5]. Even if genuine causal relationships between input and output can be established, we still do not know how the ML model produces a certain output. To answer this question, ideally, we would like to know the causal role played by (at least) some of the components making up the model. We would like to know about at least some processes, mechanisms, constraints, or structural dependencies inside of the model rather than merely probing the ML model as a black box post hoc from the outside. While XAI methods can yield information that seems plausible and physically meaningful, this information may be irrelevant with respect to how the model actually arrived at a given decision or prediction (Rudin, 2019; Baron, 2023). This, in turn, can undermine our trust in the model for future applications. In contrast, with component-level understanding, the causal knowledge is more secure and can also inform future development and improvement of the model in question and ML models in general.

## 3  Understanding and intelligibility in CMIP

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

Later, in one of the diagnostic subprojects following AMIP, Gleckler et al. (1995) attributed incorrect calculations of ocean heat transport to the models' representations of cloud radiative effects. They first found that the models' implied ocean heat transport was partially in the wrong direction – northward in the Southern Hemisphere. They inferred that cloud radiative effects were the culprit, explicitly noting that atmospheric GCMs at the time of their writing were "known to disagree considerably in their simulations of the effects of clouds on the Earth's radiation budget (Cess et al., 1989), and hence the effects of simulated cloud–radiation in-

teractions on the implied meridional energy transports [were] immediately suspect" (Gleckler et al., 1995, p. 793). They recalculated ocean heat transport using a hybrid of model data and observational data. When they did this, they fixed the error – ocean heat transport turned poleward. The observational data used to fix the error were on cloud radiative effects. In other words, they substituted the output data linked to the problematic cloud parameterizations (a component of the models) with observational data on cloud radiative effects. This substitution resulted in a better fit with observations of and physical background knowledge of ocean heat transport.

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

Many more recent cases of error diagnosis also aim to identify problematic parameterizations (see, e.g., Hall and Qu, 2006; O'Brien et al., 2013; Pitari et al., 2014; Bukovsky et al., 2017; Gettelman et al., 2019; but see Neelin et al., 2023, for current challenges). In CMIP6 in particular, there is an increased focus on process-level analysis (Eyring et al., 2019; Maloney et al., 2019). In process-level analysis, scientists examine bias in the simulation of particular processes which are, in turn, linked to one or more parameterizations, i.e., components within a whole GCM.[6] Moreover, CMIP-endorsed model intercomparison projects (MIPs) also center on particular processes or parameterizations, such as cloud feedback MIPs (Webb et al., 2017) and land surface, snow, and soil moisture MIPs (van den Hurk et al., 2016).[7]

The practice of updating model parameterizations during model development also demonstrates an interest (and success) in achieving component-level understanding. We provide two examples here: one associated with the radiative transfer parameterization in the Community Atmosphere Model and another associated with the physical representation of stratocumulus clouds in boundary layer parameterizations. With respect to the radiative transfer component (pa-

---

[5]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.
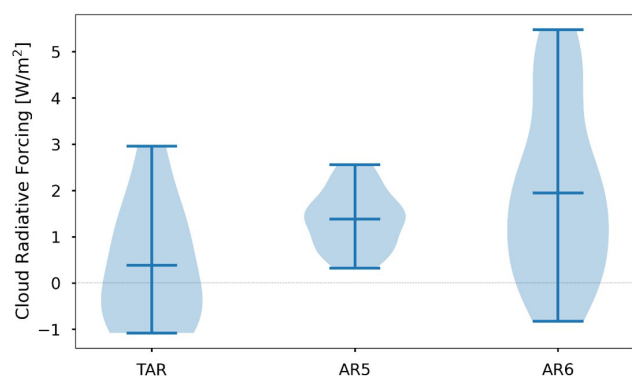
[6]Note that while processes and model components are linked, neither is reducible to the other; e.g., a coupler is a component in a GCM but it is not a real-world climate process. Conversely, there is no cloud feedback parameterization but cloud feedbacks are a real-world climate process.

[7]These examples are in stark contrast to the pessimism about understanding climate models that some philosophers of science have emphasized (Lenhard and Winsberg, 2010) and others have rebutted (Carrier and Lenhard, 2019; Easterbrook, 2023; Frigg et al., 2015; O'Loughlin, 2023; Touzé-Peiffer et al., 2020).

**Figure 3.** 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 Fig. 7.2 of Stocker et al. (2001) 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.

rameterization), Collins et al. (2002) noted that, at the time their paper was written, studies had "demonstrated that the longwave cooling rates and thermodynamic state simulated by GCMs are sensitive to the treatment of water vapor line strengths." Collins et al. used this knowledge – along with updated information about absorption and emission of thermal radiation by water vapor – to update the radiation parameterization in the Community Atmosphere Model. This component-level improvement led to substantial improvements in the models' simulated climate.

Regarding stratocumulus cloud parameterization in climate models, targeted developments following the Third Intergovernmental Panel on Climate Change (IPCC) Assessment Report reduced uncertainty in estimates of cloud feedbacks to the extent that the Sixth IPCC Assessment Report now states with high confidence that "future changes in clouds will, overall, cause additional warming" (p. 1022). This systematic change in cloud radiative forcing is demonstrated in Fig. 3. It was not clear in the Third IPCC Assessment Report (TAR) whether cloud feedbacks were positive or negative, and the TAR noted in particular that the "difficulty in simulation of observed boundary layer cloud properties is a clear testimony of the still inadequate representation of boundary-layer processes" (Stocker, 2001, p. 273). Around this time, researchers developed improved bound-

ary layer parameterizations with the goal of improving the representation of low boundary layer clouds. For instance, Grenier and Bretherton (2001) built on a standard 1.5-order boundary layer turbulence parameterization in which turbulent mixing is treated as a diffusive process related to the amount of turbulent kinetic energy (TKE) and in which TKE is treated as a conservative, prognostic quantity. Their key additions to the 1.5-order turbulence approach were (1) a more accurate numerical treatment of diffusion in the vicinity of step-function-like jumps in temperature and humidity (inversions) and (2) the contribution of cloud-top radiative cooling to the production of TKE. These two ingredients allow the turbulence parameterization to emulate the physics that drive stratocumulus clouds. Variations on the parameterization of Grenier and Bretherton (2001) and other similarly sophisticated boundary layer parameterizations have been included in numerous weather and climate models, leading to improvements in the simulation of stratocumulus clouds specifically and general improvements in model climatology.

In certain circumstances component-level responsibility for particular model biases can be determined. As an example, the Community Earth System Model 2 (CESM2) was recognized as exhibiting a climate sensitivity that is too large – one that did not appear in standard CMIP simulations. This behavior was discovered in a surprising way. Zhu et al. (2021) showed an instability in the simulation of the Last Glacial Maximum, a much colder period than the present day, using CESM2. This instability did not exist in CESM. By reverting to the original component-level microphysics scheme the model behaved as expected, and erroneous specifications of microphysical particle concentrations were discovered and remedied. More generally, the understanding and observational constraint of ice microphysics is a challenge as demonstrated by the very large variations in ice water path across CMIP models. Using perturbed parameter estimation (PPE; e.g., Eidhammer et al., 2024) can also reveal component-level sensitivities and shortcomings.

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

The climate modeling community has long strived for component-level understanding and intelligibility. This is especially evident in the work on climate model hierarchies, i.e., a group of models which spans a range of complexity and comprehensiveness (Jeevanjee et al., 2017). Writing nearly 2 decades ago, Issac Held (2005) identified model hierarchies as necessary if we wish to understand both the climate system and complex climate models.

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

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

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

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

These appeals to climate model hierarchies highlight how component-level understanding is a long-standing goal in climate modeling (see also Katzav and Parker, 2015). This is not to say that component-level understanding automatically translates to understanding all model behaviors. Emergent properties such as equilibrium climate sensitivity (ECS) may elude explanation. Even when components such as cloud parameterizations are appealed to as causally relevant for higher ECS values (e.g., Zelinka et al., 2020), it must be granted that these cloud parameterizations interact with other components and pieces of the overall GCM. That is, GCMs exhibit fuzzy modularity – sub-model behaviors do not add up linearly or in an easy-to-understand way (Lenhard and Winsberg, 2010). So, there may be a more complete explanation detailing how, as a whole, the GCM simulates a higher ECS. Producing a complete explanation may prove elusive, however, to the extent that GCMs are epistemically opaque or have such a high degree of complexity that human minds cannot track all of the relevant information (Humphreys, 2009).[8]

---

[8]This complexity includes both the impossibility of fully knowing a climate model's code in its entirety and the impossibility of being able to follow the calculations as the model steps forward in time. With today's GCMs, humans can do neither of these things.

Therefore, we do not regard our three proposed types of understanding as exhaustive – perhaps a component interaction or structural type of understanding ought to be theorized and strived for as well.

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

We end this section with a brief discussion distinguishing between component-level and statistical understanding. Overall, our analysis is in the same spirit as that of Knüsel and Baumberger (2020), who argue that data-driven models and dynamical models alike can be understood through manipulating the model so that modelers can qualitatively anticipate model behaviors. However, not all manipulations are equal. Manipulating input data and seeing associated changes in output data does not tell you how the model produces its output. The hierarchy of understanding we propose – instrumental, statistical, and component-level – concerns the degree to which and ways in which a model is intelligible or graspable (Jebeile et al., 2021; Knüsel and Baumberger, 2020). Complex models are intelligible or graspable just in the case that, and to the degree that, their behavior can be qualitatively anticipated or explained (De Regt and Dieks, 2005; Lenhard, 2006). From our perspective, component-level understanding puts scientists into a position to better anticipate and better explain model behavior. In general, statistical understanding can help us answer questions such as "do the input–output relations of the model make sense and, if so, in what way do they make sense?" This is great for finding out whether the model's behavior is consistent with expectations across a variety of cases. This may also involve manipulating input and examining associated changes in output to better anticipate future model behavior (Jebeile et al., 2021; Knüsel and Baumberger, 2020). However, this is distinct from learning why the model behaves the way it does. To answer this distinct question, we need to know how the model is working, which, in turn, involves knowing something about the pieces making up the model. Hence, component-level understanding is called for. This is exactly the type of understanding that we see aimed for, and often grasped, in CMIP experiments.

Component-level understanding often involves a different kind of knowledge related to model architecture and beyond input–output relationships. On the one hand it can demonstrate that you know what role the component is playing in the model – this shows some knowledge of model building. It may also be helpful for answering a wider range of what-if-things-had-been-different questions. Finally, and potentially the clearest benefit of component-level understanding, it can

tell one what needs to be fixed in cases of error. This should produce additional trust in the modeling enterprise more generally.[9]

## 4 Lessons learned: examples of component-level understanding in ML

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

### 4.1 Attributing model success with physics-informed machine learning

Our first example involves physics-informed machine learning, i.e., machine learning incorporated with domain knowledge and physical principles (Kashinath et al., 2021). Model success can be attributed to a specific component in a neural net if it is known that said component in the neural net is performing a physically relevant role for a given modeling task.

Beucler et al. (2019, 2021) augment a neural net's architecture via layers which enforce conservation laws that are important for emulating convection (see Fig. 2, panel a). These laws include enthalpy conservation, column-integrated water conservation, and both longwave and shortwave radiation conservation. The conservation laws are enforced "to machine precision" (Beucler et al., 2021). Following Beucler et al. (2019) and because this neural net has a physics-informed architecture, we will use the acronym NNA. NNA is trained on aqua-planet simulation data from the Super-Parameterized Community Atmosphere Model 3.0. NNA's results are compared with those of two other neural nets: one unconstrained by physics (NNU) and another "softly" constrained through a penalization term in the loss function (NNL; see Beucler et al., 2019, for further discussion).

All three NNs are evaluated based on the mean squared errors (MSEs) of their predictions and based on whether their output violates physics conservation laws (physical constraint penalty $P$, given in units of $W^2 m^{-4}$). While NNU has the highest performance in a baseline climate – i.e., a climate well-represented by the training data – NNA and

NNL each outperform NNU in a 4 K warmer climate (see Beucler et al., 2019, Table 1), which is impressive since generalizing into warmer climates is particularly challenging for ML models (Rasp et al., 2018; Li, 2023). These results may indicate that NNU performed better in the baseline climate for the "wrong" reasons. Indeed, NNU heavily violated the physical constraints in both the baseline ($P = 458 \pm 5 \times 10^2 W^2 m^{-4}$) and the 4 K warmer climate cases ($P = 3 \times 10^5 \pm 1 \times 10^6 W^2 m^{-4}$). Compare these to the NNA case (baseline: $P = 7 \times 10^{-10} \pm 1 \times 10^{-9} W^2 m^{-4}$; 4 K warmer: $P = 2 \times 10^{-9} \pm 5 \times 10^{-9} W^2 m^{-4}$).

Beucler et al. (2021) further show that NNA predicts the total thermodynamic tendency in the enthalpy conservation equation more accurately than the other NNs – "by an amount closely related to how much each NN violates enthalpy conservation" (p. 5). The particular layer in NNA responsible for enthalpy conservation is obviously the explanation for this result. This case therefore exemplifies component-level understanding, which was straightforward because of Beucler et al.'s choice of model design.

It should be noted that both NNA and NNL perform well in the 4 K warmer climate and, more generally, "[e]nforcing constraints, whether in the architecture or the loss function, can systematically reduce the error of variables that appear in the constraints" (Beucler et al., 2021, p. 5). This suggests that, when thinking purely about model performance, physical constraints do not necessarily need to be implemented in the model's architecture. However, compared with NNL, Beucler et al.'s use of NNA facilitates straightforward component-level understanding. The component-level understanding is straightforward because we know that, by virtue of the physics knowledge built into the model's architecture, NNA obeys conservation laws as it is trained and as it is tested. We can draw an analogy with dynamical climate models. NNL is to NNA as bias-corrected GCM simulations are to ones which capture relevant physical processes with high fidelity to begin with. Knowing that a model produces a physically consistent answer for physical reasons is a stronger basis for trust than merely knowing that a model produces physically consistent answers due to post hoc bias correction.

### 4.2 Explaining model error in a case of Fourier neural operators

Another example involves a recent development in using machine learning to solve partial differential equations: the Fourier neural operator (FNO) pioneered by Li et al. (2021). The innovation of FNO is the application of Fourier transforms to enable CNN-based layers that learn "solution operators" to partial differential equations in a scale-invariant way. Building on Li et al. (2021), Pathak et al. (2022) demonstrated that training an FNO network on output from a numerical weather prediction (NWP) model produced a machine learning model that emulates NWP models with high fidelity and efficiency. A key challenge noted by Pathak et

---

[9]This is not to say that component-level understanding is necessarily superior to statistical understanding. For example, 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.

al. (2022), however, was a numerical instability that limited application of the FNO model to forecasts of lengths less than 10 d.

Analysis of the instability ultimately led the group to hypothesize that the instability was due to a specific component of the FNO model: the Fourier transform itself. The problem they identified was that the sine and cosine functions employed in Fourier transforms are the eigenfunctions of the Laplace operator on a doubly periodic, Euclidean geometry, whereas the desired problem (i.e., NWP) is intrinsic to an approximately spherical geometry. In essence, the Earth's poles represent a singularity that Fourier transforms on a latitude–longitude grid are not well-equipped to handle. Bonev et al. (2023) adapt the FNO approach to spherical geometry by utilizing spherical harmonic transforms with the Laplace-operator eigenfunctions for spherical geometries as basis functions in lieu of Fourier transforms. These eigenfunctions, the spherical harmonic functions, smoothly handle the poles as a natural part of their formulation. Bonev et al. (2023) report that the application of spherical harmonic transforms, rather than Fourier transforms, results in a model that is numerically stable up to 1 year.

The application of spherical transformations stabilizes the FNO model. Bonev et al. were able to fix the FNO because they could pinpoint the Fourier transformations, a component of the FNO model, demonstrating scientists' component-level understanding.[10]

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

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

The example demonstrates component-level understanding because the units in question are manipulated. Components within the architecture of the model are turned on and off and the resultant effects are observed.[11] This puts us in a position to say, for example, that "these neurons are responsible for generating images of trees, and we know this because turning more of these neurons on yields an image with more trees (or bigger trees) and vice versa. Moreover, the other aspects of the image are unchanged no matter what

we do to these neurons". Bau et al. (2018) also show that visual artifacts are causally linked to particular units and can be removed using this causal knowledge.

This case is analogous to the study from Gleckler et al. (1995) as described in Sect. 3 above. Recall that the cloud radiative effects from the GCMs were "turned off" (substituted out and replaced with observational data) and the calculations of ocean heat transport improved. Scientists could make sense of model error because they knew that a certain deficiency in GCMs, at the time, involved components of the GCMs responsible for representing clouds. In the same way, Bau et al. (2018) are able to intervene on generations of images by linking units in their model to particular types of image classes and examining what happens to the overall image when these units are manipulated. Note that this is distinct from the closely related method of ablating specific subsets of input data, which is more closely aligned with XAI and can therefore yield statistical understanding (see, e.g., Brenowitz et al., 2020; Park et al., 2022).

While GAN dissect is not typically used in climate science research, GANs are beginning to be adopted for some climate applications (Beroche, 2021; Besombes et al., 2021). Additionally, there are potential future applications such as in atmospheric river detection (Mahesh et al., 2024). In any case, this example demonstrates yet again how component-level understanding is achievable with ML.

## 5 Discussion and recommendations for practice

In this review and perspective paper we have argued that component-level understanding ought to be strived for in ML-driven climate science. The value of component-level understanding is especially evident in the FNO problem described previously (Sect. 4.2 above). Instrumental understanding allowed the group to identify a performance issue (numerical "issues" in the polar regions) that led to numerical instability. While the group did not employ any XAI – statistical understanding – approaches, it is clear that they would have been of limited value in identifying the underlying cause of the numerical instability, since XAI methods only probe input–output mappings. Ultimately the problem was identified and later solved by applying component-level understanding of the FNO network: knowledge that a component of the network implicitly (and incorrectly) assumed a Euclidean geometry for a problem on a spherical domain.

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

– *Robustness* of the explanation is determined given small perturbations to input.

---

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

[11]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).

– *Faithfulness* is determined by comparing the predictions of perturbed input and those of unperturbed input to determine if a feature deemed important by the XAI method does in fact change the network prediction.

– *Randomization* measures how the explanation changes by perturbing the network weights; similar to the robustness metric, the thinking is that "the explanation of an input $x$ should change if the model changes or if a different class is explained" (Bommer et al., 2023, p. 8).

– *Localization* measures agreement between the explanation and a user-defined region of interest.

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

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

Therefore, our recommendation is that ML-driven climate science strive for component-level understanding. This will aid in evaluating the credibility of model results, in diagnosing model error, and in model development. The clearest path to component-level understanding in ML-driven climate science would likely involve climate scientists building, or helping build, the ML models that are used for their research and implementing physics-based and other background knowledge (Kashinath et al., 2021; Cuomo et al., 2022). Clear standards could also be developed for documenting ML architecture, training procedures, and past analyses, including error diagnoses (O'Loughlin, 2023). Perhaps a model intercomparison project could be developed to systematically evaluate ML behavior across diverse groups of researchers. Lastly, with component-level understanding as a goal to strive for, scientists can better develop hybrid models where both ML and dynamic modeling components are employed.

An increasing range of free or low-cost, high-quality resources are now available to enable researchers who are not (yet) experts in ML to gain a deep and practical level of understanding of modern ML model designs and applications. Some examples of free, high-quality resources include the following.[12]

– Practical Deep Learning for Coders – 1: Getting started (fast.ai): https://course.fast.ai/Lessons/lesson1.html

[12] All links last accessed 16 December 2024.

– Related: GitHub – fastai/fastbook: The fastai book, published as Jupyter Notebooks: https://github.com/fastai/fastbook

– Introduction – Hugging Face NLP Course: https://huggingface.co/learn/nlp-course/chapter1/1

– How Diffusion Models Work – DeepLearning.AI: https://www.deeplearning.ai/short-courses/how-diffusion-models-work/

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

# References

Balmaceda-Huarte, R., Baño-Medina, J., Olmo, M. E., and Bettolli, M. L.: On the use of convolutional neural networks for downscaling daily temperatures over southern South America in a climate change scenario, Clim. Dynam., 62, 383–397, https://doi.org/10.1007/s00382-023-06912-6, 2023.

Barnes, E. A., Barnes, R. J., Martin, Z. K., and Rader, J. K.: This Looks Like That There: Interpretable Neural Networks for Image Tasks When Location Matters, Artif. Intell. Earth Syst., 1, e220001, https://doi.org/10.1175/AIES-D-22-0001.1, 2022.

Baron, S.: Explainable AI and Causal Understanding: Counterfactual Approaches Considered, Minds Mach., 33, 347–377, https://doi.org/10.1007/s11023-023-09637-x, 2023.

Bau, D., Zhu, J.-Y., Strobelt, H., Zhou, B., Tenenbaum, J. B., Freeman, W. T., and Torralba, A.: GAN Dissection: Visualizing and Understanding Generative Adversarial Networks, arXiv [preprint], https://doi.org/10.48550/arXiv.1811.10597, 8 December 2018.

Baumberger, C., Knutti, R., and Hadorn, G. H.: Building confidence in climate model projections: an analysis of inferences from fit, WIREs Clim. Change, 8, e454, https://doi.org/10.1002/wcc.454, 2017.

Beroche, H.: Generative Adversarial Networks for Climate Change Scenarios, URBAN AI, https://urbanai.fr/generative-adversarial-networks-for-climate-change-scenarios/ (last access: 16 December 2024), 2021.

Besombes, C., Pannekoucke, O., Lapeyre, C., Sanderson, B., and Thual, O.: Producing realistic climate data with generative adversarial networks, Nonlin. Processes Geophys., 28, 347–370, https://doi.org/10.5194/npg-28-347-2021, 2021.

Beucler, T., Rasp, S., Pritchard, M., and Gentine, P.: Achieving Conservation of Energy in Neural Network Emulators for Climate Modeling, arXiv [preprint], https://doi.org/10.48550/arXiv.1906.06622, 15 June 2019.

Beucler, T., Pritchard, M., Rasp, S., Ott, J., Baldi, P., and Gentine, P.: Enforcing Analytic Constraints in Neural Networks Emulating Physical Systems, Phys. Rev. Lett., 126, 098302, https://doi.org/10.1103/PhysRevLett.126.098302, 2021.

Bommer, P., Kretschmer, M., Hedström, A., Bareeva, D., and Höhne, M. M.-C.: Finding the right XAI method – A Guide for the Evaluation and Ranking of Explainable AI Methods in Climate Science, arXiv [preprint], https://doi.org/10.48550/arXiv.2303.00652, 1 March 2023.

Bonev, B., Kurth, T., Hundt, C., Pathak, J., Baust, M., Kashinath, K., and Anandkumar, A.: Spherical Fourier Neural Operators: Learning Stable Dynamics on the Sphere, arXiv [preprint], https://doi.org/10.48550/arXiv.2306.03838, 6 June 2023.

Brenowitz, N. D., Beucler, T., Pritchard, M., and Bretherton, C. S.: Interpreting and Stabilizing Machine-Learning Parametrizations of Convection, J. Atmospheric Sci., 77, 4357–4375, https://doi.org/10.1175/JAS-D-20-0082.1, 2020.

Bukovsky, M. S., McCrary, R. R., Seth, A., and Mearns, L. O.: A Mechanistically Credible, Poleward Shift in Warm-Season Precipitation Projected for the U.S. Southern Great Plains?, J. Climate, 30, 8275–8298, https://doi.org/10.1175/JCLI-D-16-0316.1, 2017.

Carrier, M. and Lenhard, J.: Climate Models: How to Assess Their Reliability, Int. Stud. Philos. Sci., 32, 81–100, https://doi.org/10.1080/02698595.2019.1644722, 2019.

Cess, R. D., Potter, G. L., Blanchet, J. P., Boer, G. J., Ghan, S. J., Kiehl, J. T., Treut, H. L., Li, Z.-X., Liang, X.-Z., Mitchell, J. F. B., Morcrette, J.-J., Randall, D. A., Riches, M. R., Roeckner, E., Schlese, U., Slingo, A., Taylor, K. E., Washington, W. M., Wetherald, R. T., and Yagai, I.: Interpretation of Cloud-Climate Feedback as Produced by 14 Atmospheric General Circulation Models, Science, 245, 513–516, https://doi.org/10.1126/science.245.4917.513, 1989.

Chakraborty, D., Başağaoğlu, H., Gutierrez, L., and Mirchi, A.: Explainable AI reveals new hydroclimatic insights for ecosystem-centric groundwater management, Environ. Res. Lett., 16, 114024, https://doi.org/10.1088/1748-9326/ac2fde, 2021.

Cilli, R., Elia, M., D'Este, M., Giannico, V., Amoroso, N., Lombardi, A., Pantaleo, E., Monaco, A., Sanesi, G., Tangaro, S., Bellotti, R., and Lafortezza, R.: Explainable artificial intelligence (XAI) detects wildfire occurrence in the Mediterranean countries of Southern Europe, Sci. Rep., 12, 16349, https://doi.org/10.1038/s41598-022-20347-9, 2022.

Clare, M. C. A., Sonnewald, M., Lguensat, R., Deshayes, J., and Balaji, V.: Explainable Artificial Intelligence for Bayesian Neural Networks: Toward Trustworthy Predictions of Ocean Dynamics, J. Adv. Model. Earth Sy., 14, e2022MS003162, https://doi.org/10.1029/2022MS003162, 2022.

Collins, W. D., Hackney, J. K., and Edwards, D. P.: An updated parameterization for infrared emission and absorption by water vapor in the National Center for Atmospheric Research Community Atmosphere Model, J. Geophys. Res.-Atmos., 107, ACL 17-1–ACL 17-20, https://doi.org/10.1029/2001JD001365, 2002.

Cuomo, S., Di Cola, V. S., Giampaolo, F., Rozza, G., Raissi, M., and Piccialli, F.: Scientific Machine Learning Through Physics–Informed Neural Networks: Where we are and What's Next,

J. Sci. Comput., 92, 88, https://doi.org/10.1007/s10915-022-01939-z, 2022.

de Burgh-Day, C. O. and Leeuwenburg, T.: Machine learning for numerical weather and climate modelling: a review, Geosci. Model Dev., 16, 6433–6477, https://doi.org/10.5194/gmd-16-6433-2023, 2023.

De Regt, H. W.: Understanding Scientific Understanding, Oxford University Press, 321 pp., ISBN: 9780190652913, 2017.

De Regt, H. W. and Dieks, D.: A Contextual Approach to Scientific Understanding, Synthese, 144, 137–170, https://doi.org/10.1007/s11229-005-5000-4, 2005.

Diffenbaugh, N. S. and Barnes, E. A.: Data-Driven Predictions of the Time Remaining until Critical Global Warming Thresholds Are Reached, P. Natl. Acad. Sci. USA, 120, e2207183120, https://doi.org/10.1073/pnas.2207183120, 2023.

Easterbrook, S. M.: Computing the Climate: How We Know What We Know About Climate Change, Cambridge University Press, Cambridge, https://doi.org/10.1017/9781316459768, 2023.

Eidhammer, T., Gettelman, A., Thayer-Calder, K., Watson-Parris, D., Elsaesser, G., Morrison, H., van Lier-Walqui, M., Song, C., and McCoy, D.: An extensible perturbed parameter ensemble for the Community Atmosphere Model version 6, Geosci. Model Dev., 17, 7835–7853, https://doi.org/10.5194/gmd-17-7835-2024, 2024.

Eyring, V., Cox, P. M., Flato, G. M., Gleckler, P. J., Abramowitz, G., Caldwell, P., Collins, W. D., Gier, B. K., Hall, A. D., Hoffman, F. M., Hurtt, G. C., Jahn, A., Jones, C. D., Klein, S. A., Krasting, J. P., Kwiatkowski, L., Lorenz, R., Maloney, E., Meehl, G. A., Pendergrass, A. G., Pincus, R., Ruane, A. C., Russell, J. L., Sanderson, B. M., Santer, B. D., Sherwood, S. C., Simpson, I. R., Stouffer, R. J., and Williamson, M. S.: Taking climate model evaluation to the next level, Nat. Clim. Change, 9, 102–110, https://doi.org/10.1038/s41558-018-0355-y, 2019.

Felsche, E. and Ludwig, R.: Applying machine learning for drought prediction in a perfect model framework using data from a large ensemble of climate simulations, Nat. Hazards Earth Syst. Sci., 21, 3679–3691, https://doi.org/10.5194/nhess-21-3679-2021, 2021.

Forster, P., Storelvmo, T., Armour, K., Collins, W., Dufresne, J.-L., Frame, D., Lunt, D. J., Mauritsen, T., Palmer, M. D., Watanabe, M., Wild, M., and Zhang, H.: The Earth's Energy Budget, Climate Feedbacks, and Climate Sensitivity, in: Climate Change 2021: The Physical Science Basis. Contribution of Working Group I to the Sixth Assessment Report of the Intergovernmental Panel on Climate Change, edited by: Masson-Delmotte, V., Zhai, P., Pirani, A., Connors, S. L., Péan, C., Berger, S., Caud, N., Chen, Y., Goldfarb, L., Gomis, M. I., Huang, M., Leitzell, K., Lonnoy, E., Matthews, J. B. R., Maycock, T. K., Waterfield, T., Yelekçi, O., Yu, R., and Zhou, B., Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA, 923–1054, https://doi.org/10.1017/9781009157896.009, 2021.

Frigg, R., Thompson, E., and Werndl, C.: Philosophy of Climate Science Part II: Modelling Climate Change, Philos. Compass, 10, 965–977, https://doi.org/10.1111/phc3.12297, 2015.

Gates, W. L.: AMIP: The Atmospheric Model Intercomparison Project, B. Am. Meteorol. Soc., 73, 1962–1970, https://doi.org/10.1175/1520-0477(1992)073<1962:ATAMIP>2.0.CO;2, 1992.

Gettelman, A., Hannay, C., Bacmeister, J. T., Neale, R. B., Pendergrass, A. G., Danabasoglu, G., Lamarque, J.-F., Fasullo, J. T., Bailey, D. A., Lawrence, D. M., and Mills, M. J.: High Climate Sensitivity in the Community Earth System Model Version 2 (CESM2), Geophys. Res. Lett., 46, 8329–8337, https://doi.org/10.1029/2019GL083978, 2019.

Gleckler, P. J., Randall, D. A., Boer, G., Colman, R., Dix, M., Galin, V., Helfand, M., Kiehl, J., Kitoh, A., Lau, W., Liang, X.-Y., Lykossov, V., McAvaney, B., Miyakoda, K., Planton, S., and Stern, W.: Cloud-radiative effects on implied oceanic energy transports as simulated by Atmospheric General Circulation Models, Geophys. Res. Lett., 22, 791–794, https://doi.org/10.1029/95GL00113, 1995.

Gleckler, P. J., Taylor, K. E., and Doutriaux, C.: Performance metrics for climate models, J. Geophys. Res.-Atmos., 113, D06104, https://doi.org/10.1029/2007JD008972, 2008.

González-Abad, J., Baño-Medina, J., and Gutiérrez, J. M.: Using Explainability to Inform Statistical Downscaling Based on Deep Learning Beyond Standard Validation Approaches, arXiv [preprint], https://doi.org/10.48550/arXiv.2302.01771, 3 February 2023.

Gordon, E. M., Barnes, E. A., and Hurrell, J. W.: Oceanic Harbingers of Pacific Decadal Oscillation Predictability in CESM2 Detected by Neural Networks, Geophys. Res. Lett., 48, e2021GL095392, https://doi.org/10.1029/2021GL095392, 2021.

Grenier, H. and Bretherton, C. S.: A Moist PBL Parameterization for Large-Scale Models and Its Application to Subtropical Cloud-Topped Marine Boundary Layers, Mon. Weather Rev., 129, 357–377, https://doi.org/10.1175/1520-0493(2001)129<0357:AMPPFL>2.0.CO;2, 2001.

Grundner, A., Beucler, T., Gentine, P., Iglesias-Suarez, F., Giorgetta, M. A., and Eyring, V.: Deep Learning Based Cloud Cover Parameterization for ICON, J. Adv. Model. Earth Sy., 14, e2021MS002959, https://doi.org/10.1029/2021MS002959, 2022.

Hall, A. and Qu, X.: Using the current seasonal cycle to constrain snow albedo feedback in future climate change, Geophys. Res. Lett., 33, L03502, https://doi.org/10.1029/2005GL025127, 2006.

Ham, Y.-G., Kim, J.-H., and Luo, J.-J.: Deep learning for multi-year ENSO forecasts, Nature, 573, 568–572, https://doi.org/10.1038/s41586-019-1559-7, 2019.

Hausfather, Z., Drake, H. F., Abbott, T., and Schmidt, G. A.: Evaluating the Performance of Past Climate Model Projections, Geophys. Res. Lett., 47, e2019GL085378, https://doi.org/10.1029/2019GL085378, 2020.

He, R., Zhang, L., and Chew, A. W. Z.: Data-driven multi-step prediction and analysis of monthly rainfall using explainable deep learning, Expert Syst. Appl., 235, 121160, https://doi.org/10.1016/j.eswa.2023.121160, 2024.

Hedström, A., Weber, L., Krakowczyk, D., Bareeva, D., Motzkus, F., Samek, W., Lapuschkin, S., and Höhne, M. M.-C.: Quantus: An explainable ai toolkit for responsible evaluation of neural network explanations and beyond, J. Mach. Learn. Res., 24, 1–11, http://jmlr.org/papers/v24/22-0142.html (last access: 16 December 2024), 2023.

Hegerl, G. C., Zwiers, F. W., Braconnot, P., Gillett, N. P., Luo, Y., Orsini, J. A. M., Nicholls, N., Penner, J. E., Stott, P. A., Allen, M., Ammann, C., Andronova, N., Betts, R. A., Clement, A., Collins, W. D., Crooks, S., Delworth, T. L., Forest, C., Forster,

P., Goosse, H., Gregory, J. M., Harvey, D., Jones, G. S., Joos, F., Kenyon, J., Kettleborough, J., Kharin, V., Knutti, R., Lambert, F. H., Lavine, M., Lee, T. C. K., Levinson, D., Masson-Delmotte, V., Nozawa, T., Otto-Bliesner, B., Pierce, D., Power, S., Rind, D., Rotstayn, L., Santer, B. D., Senior, C., Sexton, D., Stark, S., Stone, D. A., Tett, S., Thorne, P., van Dorland, R., Wong, T., Xu, L., Zhang, X., Zorita, E., Karoly, D. J., Ogallo, L., and Planton, S.: Understanding and Attributing Climate Change, in: Climate Change 2007: The Physical Science Basis. Contribution of Working Group I to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change, edited by: Solomon, S., Qin, D., Manning, M., Chen, Z., Marquis, M., Avery, K. B., Tignor, M., and Miller, H. L., Cambridge University Press, 84, https://www.ipcc.ch/report/ar4/wg1/ (last access: 16 December 2024), 2007.

Held, I. M.: The Gap between Simulation and Understanding in Climate Modeling, B. Am. Meteorol. Soc., 86, 1609–1614, https://doi.org/10.1175/BAMS-86-11-1609, 2005.

Hourdin, F., Grandpeix, J.-Y., Rio, C., Bony, S., Jam, A., Cheruy, F., Rochetin, N., Fairhead, L., Idelkadi, A., Musat, I., Dufresne, J.-L., Lahellec, A., Lefebvre, M.-P., and Roehrig, R.: LMDZ5B: the atmospheric component of the IPSL climate model with revisited parameterizations for clouds and convection, Clim. Dynam., 40, 2193–2222, https://doi.org/10.1007/s00382-012-1343-y, 2013.

Humphreys, P.: The philosophical novelty of computer simulation methods, Synthese, 169, 615–626, https://doi.org/10.1007/s11229-008-9435-2, 2009.

Jebeile, J., Lam, V., and Räz, T.: Understanding climate change with statistical downscaling and machine learning, Synthese, 199, 1877–1897, https://doi.org/10.1007/s11229-020-02865-z, 2021.

Jeevanjee, N., Hassanzadeh, P., Hill, S., and Sheshadri, A.: A perspective on climate model hierarchies, J. Adv. Model. Earth Sy., 9, 1760–1771, https://doi.org/10.1002/2017MS001038, 2017.

Kashinath, K., Mustafa, M., Albert, A., Wu, J.-L., Jiang, C., Esmaeilzadeh, S., Azizzadenesheli, K., Wang, R., Chattopadhyay, A., Singh, A., Manepalli, A., Chirila, D., Yu, R., Walters, R., White, B., Xiao, H., Tchelepi, H. A., Marcus, P., Anandkumar, A., Hassanzadeh, P., and Prabhat, null: Physics-informed machine learning: case studies for weather and climate modelling, Philos. T. R. Soc. A, 379, 20200093, https://doi.org/10.1098/rsta.2020.0093, 2021.

Katzav, J. and Parker, W. S.: The future of climate modeling, Climatic Change, 132, 475–487, https://doi.org/10.1007/s10584-015-1435-x, 2015.

Knüsel, B. and Baumberger, C.: Understanding climate phenomena with data-driven models, Stud. Hist. Philos. Sci. Part A, 84, 46–56, https://doi.org/10.1016/j.shpsa.2020.08.003, 2020.

Knutti, R.: Why are climate models reproducing the observed global surface warming so well?, Geophys. Res. Lett., 35, L18704, https://doi.org/10.1029/2008GL034932, 2008.

Kravitz, B., Robock, A., Forster, P. M., Haywood, J. M., Lawrence, M. G., and Schmidt, H.: An overview of the Geoengineering Model Intercomparison Project (GeoMIP), J. Geophys. Res.-Atmos., 118, 13103–13107, https://doi.org/10.1002/2013JD020569, 2013.

Labe, Z. M. and Barnes, E. A.: Comparison of Climate Model Large Ensembles With Observations in the Arctic Using Simple Neural Networks, Earth Space Sci., 9, e2022EA002348, https://doi.org/10.1029/2022EA002348, 2022a.

Labe, Z. M. and Barnes, E. A.: Predicting Slowdowns in Decadal Climate Warming Trends With Explainable Neural Networks, Geophys. Res. Lett., 49, e2022GL098173, https://doi.org/10.1029/2022GL098173, 2022b.

Krishna, S., Han, T., Gu, A., Pombra, J., Jabbari, S., Wu, S., and Lakkaraju, H.: The Disagreement Problem in Explainable Machine Learning: A Practitioner's Perspective, arXiv [preprint], https://doi.org/10.48550/arXiv.2202.01602, 8 February 2022.

Labe, Z. M. and Barnes, E. A.: Detecting Climate Signals Using Explainable AI With Single-Forcing Large Ensembles, J. Adv. Model. Earth Sy., 13, e2021MS002464, https://doi.org/10.1029/2021MS002464, 2021.

Lenhard, J.: Surprised by a Nanowire: Simulation, Control, and Understanding, Philos. Sci., 73, 605–616, https://doi.org/10.1086/518330, 2006.

Lenhard, J. and Winsberg, E.: Holism, entrenchment, and the future of climate model pluralism, Stud. Hist. Philos. Sci. Part B Stud. Hist. Philos. Mod. Phys., 41, 253–262, https://doi.org/10.1016/j.shpsb.2010.07.001, 2010.

Li, D.: Machines Learn Better with Better Data Ontology: Lessons from Philosophy of Induction and Machine Learning Practice, Minds Mach., 33, 429–450, https://doi.org/10.1007/s11023-023-09639-9, 2023.

Li, W., Migliavacca, M., Forkel, M., Denissen, J. M. C., Reichstein, M., Yang, H., Duveiller, G., Weber, U., and Orth, R.: Widespread increasing vegetation sensitivity to soil moisture, Nat. Commun., 13, 3959, https://doi.org/10.1038/s41467-022-31667-9, 2022.

Li, Z., Kovachki, N., Azizzadenesheli, K., Liu, B., Bhattacharya, K., Stuart, A., and Anandkumar, A.: Fourier Neural Operator for Parametric Partial Differential Equations, arXiv [preprint], https://doi.org/10.48550/arXiv.2010.08895, 16 May 2021.

Lin, Q.-J., Mayta, V. C., and Adames Corraliza, Á. F.: Assessment of the Madden-Julian Oscillation in CMIP6 Models Based on Moisture Mode Theory, Geophys. Res. Lett., 51, e2023GL106693, https://doi.org/10.1029/2023GL106693, 2024.

Lipton, Z. C.: The mythos of model interpretability (2016), arXiv [preprint], https://doi.org/10.48550/arXiv.1606.03490, 2016.

Liu, Y., Duffy, K., Dy, J. G., and Ganguly, A. R.: Explainable deep learning for insights in El Niño and river flows, Nat. Commun., 14, 339, https://doi.org/10.1038/s41467-023-35968-5, 2023.

Lloyd, E. A.: Confirmation and Robustness of Climate Models, Philos. Sci., 77, 971–984, https://doi.org/10.1086/657427, 2010.

Lloyd, E. A.: Model robustness as a confirmatory virtue: The case of climate science, Stud. Hist. Philos. Sci. Part A, 49, 58–68, https://doi.org/10.1016/j.shpsa.2014.12.002, 2015.

Mahesh, A., O'Brien, T. A., Loring, B., Elbashandy, A., Boos, W., and Collins, W. D.: Identifying atmospheric rivers and their poleward latent heat transport with generalizable neural networks: ARCNNv1, Geosci. Model Dev., 17, 3533–3557, https://doi.org/10.5194/gmd-17-3533-2024, 2024.

Maloney, E. D., Gettelman, A., Ming, Y., Neelin, J. D., Barrie, D., Mariotti, A., Chen, C.-C., Coleman, D. R. B., Kuo, Y.-H., Singh, B., Annamalai, H., Berg, A., Booth, J. F., Camargo, S. J., Dai, A., Gonzalez, A., Hafner, J., Jiang, X., Jing, X., Kim, D., Kumar, A., Moon, Y., Naud, C. M., Sobel, A. H., Suzuki, K., Wang, F., Wang, J., Wing, A. A., Xu, X., and Zhao, M.: Process-Oriented Evaluation of Climate and Weather Forecasting Models, B. Am. Meteorol. Soc., 100, 1665–1686, https://doi.org/10.1175/BAMS-D-18-0042.1, 2019.

Mamalakis, A., Ebert-Uphoff, I., and Barnes, E. A.: Explainable Artificial Intelligence in Meteorology and Climate Science: Model Fine-Tuning, Calibrating Trust and Learning New Science, in: xxAI – Beyond Explainable AI: International Workshop, Held in Conjunction with ICML 2020, July 18, 2020, Vienna, Austria, Revised and Extended Papers, edited by: Holzinger, A., Goebel, R., Fong, R., Moon, T., Müller, K.-R., and Samek, W., Springer International Publishing, Cham, 315–339, https://doi.org/10.1007/978-3-031-04083-2_16, 2022a.

Mamalakis, A., Ebert-Uphoff, I., and Barnes, E. A.: Neural network attribution methods for problems in geoscience: A novel synthetic benchmark dataset, Environ. Data Sci., 1, e8, https://doi.org/10.1017/eds.2022.7, 2022b.

Mayernik, M. S.: Credibility via Coupling: Institutions and Infrastructures in Climate Model Intercomparisons:, Engag. Sci. Technol. Soc., 7, 10–32, https://doi.org/10.17351/ests2021.769, 2021.

McGinnis, S., Korytina, D., Bukovsky, M., McCrary, R., and Mearns, L.: Credibility Evaluation of a Convolutional Neural Net for Downscaling GCM Output over the Southern Great Plains, 2021, GC42A-03, https://ui.adsabs.harvard.edu/abs/2021AGUFMGC42A..03M%2F/abstract (last access: 3 February 2025), 2021.

Molina, M. J., O'Brien, T. A., Anderson, G., Ashfaq, M., Bennett, K. E., Collins, W. D., Dagon, K., Restrepo, J. M., and Ullrich, P. A.: A Review of Recent and Emerging Machine Learning Applications for Climate Variability and Weather Phenomena, Artif. Intell. Earth Syst., 2, 220086, https://doi.org/10.1175/AIES-D-22-0086.1, 2023.

Neelin, J. D., Krasting, J. P., Radhakrishnan, A., Liptak, J., Jackson, T., Ming, Y., Dong, W., Gettelman, A., Coleman, D. R., Maloney, E. D., Wing, A. A., Kuo, Y.-H., Ahmed, F., Ullrich, P., Bitz, C. M., Neale, R. B., Ordonez, A., and Maroon, E. A.: Process-Oriented Diagnostics: Principles, Practice, Community Development, and Common Standards, B. Am. Meteorol. Soc., 104, E1452–E1468, https://doi.org/10.1175/BAMS-D-21-0268.1, 2023.

Notz, D., Haumann, F. A., Haak, H., Jungclaus, J. H., and Marotzke, J.: Arctic sea-ice evolution as modeled by Max Planck Institute for Meteorology's Earth system model, J. Adv. Model. Earth Sy., 5, 173–194, https://doi.org/10.1002/jame.20016, 2013.

NSF AI Institute for Research on Trustworthy AI in Weather, Climate, and Coastal Oceanography (AI2ES), https://www.ai2es.org/, last access: 13 August 2024.

O'Brien, T. A., Li, F., Collins, W. D., Rauscher, S. A., Ringler, T. D., Taylor, M., Hagos, S. M., and Leung, L. R.: Observed Scaling in Clouds and Precipitation and Scale Incognizance in Regional to Global Atmospheric Models, J. Climate, 26, 9313–9333, https://doi.org/10.1175/JCLI-D-13-00005.1, 2013.

O'Loughlin, R.: Robustness reasoning in climate model comparisons, Stud. Hist. Philos. Sci. Part A, 85, 34–43, https://doi.org/10.1016/j.shpsa.2020.12.005, 2021.

O'Loughlin, R.: Diagnosing errors in climate model intercomparisons, Eur. J. Philos. Sci., 13, 20, https://doi.org/10.1007/s13194-023-00522-z, 2023.

Oreopoulos, L., Mlawer, E., Delamere, J., Shippert, T., Cole, J., Fomin, B., Iacono, M., Jin, Z., Li, J., Manners, J., Räisänen, P., Rose, F., Zhang, Y., Wilson, M. J., and Rossow, W. B.: The Continual Intercomparison of Radiation Codes: Re-

sults from Phase I, J. Geophys. Res.-Atmos., 117, D06118, https://doi.org/10.1029/2011JD016821, 2012.

Park, M., Tran, D. Q., Bak, J., and Park, S.: Advanced wildfire detection using generative adversarial network-based augmented datasets and weakly supervised object localization, Int. J. Appl. Earth Obs., 114, 103052, https://doi.org/10.1016/j.jag.2022.103052, 2022.

Pathak, J., Subramanian, S., Harrington, P., Raja, S., Chattopadhyay, A., Mardani, M., Kurth, T., Hall, D., Li, Z., Azizzadenesheli, K., Hassanzadeh, P., Kashinath, K., and Anandkumar, A.: FourCastNet: A Global Data-driven High-resolution Weather Model using Adaptive Fourier Neural Operators, arXiv [preprint], https://doi.org/10.48550/arXiv.2202.11214, 22 February 2022.

Pincus, R., Forster, P. M., and Stevens, B.: The Radiative Forcing Model Intercomparison Project (RFMIP): experimental protocol for CMIP6, Geosci. Model Dev., 9, 3447–3460, https://doi.org/10.5194/gmd-9-3447-2016, 2016.

Pitari, G., Aquila, V., Kravitz, B., Robock, A., Watanabe, S., Cionni, I., Luca, N. D., Genova, G. D., Mancini, E., and Tilmes, S.: Stratospheric ozone response to sulfate geoengineering: Results from the Geoengineering Model Intercomparison Project (GeoMIP), J. Geophys. Res.-Atmos., 119, 2629–2653, https://doi.org/10.1002/2013JD020566, 2014.

Rader, J. K., Barnes, E. A., Ebert-Uphoff, I., and Anderson, C.: Detection of Forced Change Within Combined Climate Fields Using Explainable Neural Networks, J. Adv. Model. Earth Sy., 14, e2021MS002941, https://doi.org/10.1029/2021MS002941, 2022.

Rasp, S., Pritchard, M. S., and Gentine, P.: Deep learning to represent subgrid processes in climate models, P. Natl. Acad. Sci. USA, 115, 9684–9689, https://doi.org/10.1073/pnas.1810286115, 2018.

Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., and Prabhat: Deep learning and process understanding for data-driven Earth system science, Nature, 566, 195–204, https://doi.org/10.1038/s41586-019-0912-1, 2019.

Rudin, C.: Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead, Nat. Mach. Intell., 1, 206–215, https://doi.org/10.1038/s42256-019-0048-x, 2019.

Schmidt, G. A. and Sherwood, S.: A practical philosophy of complex climate modelling, Eur. J. Philos. Sci., 5, 149–169, https://doi.org/10.1007/s13194-014-0102-9, 2015.

Schneider, S. H.: On the Carbon Dioxide–Climate Confusion, J. Atmos. Sci., 32, 2060–2066, https://doi.org/10.1175/1520-0469(1975)032<2060:OTCDC>2.0.CO;2, 1975.

Schneider, S. H.: Verification of Parameterizations in Climate Modeling, in: Report of the Study Conference on Climate Models: Performance, Intercomparison and Sensitivity Studies, edited by: Gates, W. L., World Meteorological Organization, Global Atmospheric Research Program, GARP Publications Series no. 22, Vol. 2, 728–751, 1979.

Stocker, T. F.: Working Group 1: 7.2.2.3 Boundary-layer mixing and cloudiness from the IPCC's Third Assessment Report (TAR), https://archive.ipcc.ch/ipccreports/tar/wg1/273.htm (last access: 26 October 2023), 2001.

Stocker, T. F., et al.: Physical climate processes and feedbacks, IPCC, Climate change 2001: The scientific basis. Contribution

of working group I to the third assessment report of the intergovernmental panel on climate change, Cambridge University Press, 417–470, 2001.

Subel, A., Guan, Y., Chattopadhyay, A., and Hassanzadeh, P.: Explaining the physics of transfer learning in data-driven turbulence modeling, PNAS Nexus, 2, pgad015, https://doi.org/10.1093/pnasnexus/pgad015, 2023.

Toms, B. A., Barnes, E. A., and Hurrell, J. W.: Assessing Decadal Predictability in an Earth-System Model Using Explainable Neural Networks, Geophys. Res. Lett., 48, e2021GL093842, https://doi.org/10.1029/2021GL093842, 2021.

Touzé-Peiffer, L., Barberousse, A., and Treut, H. L.: The Coupled Model Intercomparison Project: History, uses, and structural effects on climate research, WIREs Clim. Change, 11, e648, https://doi.org/10.1002/wcc.648, 2020.

van den Hurk, B., Kim, H., Krinner, G., Seneviratne, S. I., Derksen, C., Oki, T., Douville, H., Colin, J., Ducharne, A., Cheruy, F., Viovy, N., Puma, M. J., Wada, Y., Li, W., Jia, B., Alessandri, A., Lawrence, D. M., Weedon, G. P., Ellis, R., Hagemann, S., Mao, J., Flanner, M. G., Zampieri, M., Materia, S., Law, R. M., and Sheffield, J.: LS3MIP (v1.0) contribution to CMIP6: the Land Surface, Snow and Soil moisture Model Intercomparison Project – aims, setup and expected outcome, Geosci. Model Dev., 9, 2809–2832, https://doi.org/10.5194/gmd-9-2809-2016, 2016.

Wang, S., Sankaran, S., and Perdikaris, P.: Respecting causality is all you need for training physics-informed neural networks, arXiv [preprint], https://doi.org/10.48550/arXiv.2203.07404, 14 March 2022.

Webb, M. J., Andrews, T., Bodas-Salcedo, A., Bony, S., Bretherton, C. S., Chadwick, R., Chepfer, H., Douville, H., Good, P., Kay, J. E., Klein, S. A., Marchand, R., Medeiros, B., Siebesma, A. P., Skinner, C. B., Stevens, B., Tselioudis, G., Tsushima, Y., and Watanabe, M.: The Cloud Feedback Model Intercomparison Project (CFMIP) contribution to CMIP6, Geosci. Model Dev., 10, 359–384, https://doi.org/10.5194/gmd-10-359-2017, 2017.

Xue, P., Wagh, A., Ma, G., Wang, Y., Yang, Y., Liu, T., and Huang, C.: Integrating Deep Learning and Hydrodynamic Modeling to Improve the Great Lakes Forecast, Remote Sens., 14, 2640, https://doi.org/10.3390/rs14112640, 2022.

Yuan, H., Yu, H., Gui, S., and Ji, S.: Explainability in Graph Neural Networks: A Taxonomic Survey, IEEE Trans. Pattern Anal. Mach. Intell., 45, 5782–5799, https://doi.org/10.1109/TPAMI.2022.3204236, 2023.

Zelinka, M. D., Myers, T. A., McCoy, D. T., Po-Chedley, S., Caldwell, P. M., Ceppi, P., Klein, S. A., and Taylor, K. E.: Causes of Higher Climate Sensitivity in CMIP6 Models, Geophys. Res. Lett., 47, e2019GL085782, https://doi.org/10.1029/2019GL085782, 2020.

Zhu, J., Otto-Bliesner, B. L., Brady, E. C., Poulsen, C. J., Tierney, J. E., Lofverstrom, M., and DiNezio, P.: Assessment of Equilibrium Climate Sensitivity of the Community Earth System Model Version 2 Through Simulation of the Last Glacial Maximum, Geophys. Res. Lett., 48, e2020GL091220, https://doi.org/10.1029/2020GL091220, 2021.