

# 1 **Developing Guidelines for Working with Multi-Model Ensembles in** 2 **CMIP**

3 Anja Katzenberger<sup>1,2</sup>, Jhayron S. Perez-Carrasquilla<sup>3</sup>, Keighan Gemmell<sup>4</sup>, Evgenia Galytska<sup>5,6</sup>, Christine  
4 Leclerc<sup>7</sup>, Punya P<sup>8</sup>, Indrani Roy<sup>9</sup>, Arianna Varuolo-Clarke<sup>10,11</sup>, Milica Tošić<sup>12</sup>, Nina Črnivec<sup>13</sup>

5  
6 <sup>1</sup> Potsdam Institute for Climate Impact Research, Potsdam, 14473, Germany

7 <sup>2</sup> Institute of Physics and Astronomy, Potsdam University, Potsdam, 14469, Germany

8 <sup>3</sup> Atmospheric and Oceanic Science Department, University of Maryland, College Park, 20740, United States

9 <sup>4</sup> Department of Chemistry, The University of British Columbia, Vancouver, V6T 1Z4, Canada

10 <sup>5</sup> University of Bremen, Institute of Environmental Physics, Bremen, Germany

11 <sup>6</sup> Deutsches Zentrum für Luft- und Raumfahrt (DLR), Institut für Physik der Atmosphäre, Oberpfaffenhofen, Germany

12 <sup>7</sup> Department of Geography, Simon Fraser University, Burnaby, V5A 1S6, Canada

13 <sup>8</sup> Department of Earth and Space Sciences, Indian Institute of Space Science and Technology, Trivandrum, 695547, India

14 <sup>9</sup> University College London (UCL), Earth Science Department, Gower Street, London, WC1E 6BT, UK

15 <sup>10</sup> Cooperative Programs for the Advancement of Earth System Science, University Corporation for Atmospheric Research, Boulder, CO

16 <sup>11</sup> Cooperative Institute for Research in Environmental Sciences, University of Colorado, Boulder, CO

17 <sup>12</sup> Faculty of Physics, University of Belgrade, Belgrade, 11000, Serbia

18 <sup>13</sup> Faculty of Mathematics and Physics, University of Ljubljana, Ljubljana, 1000, Slovenia

19  
20 *Correspondence to:* Anja Katzenberger (anja.katzenberger@gmx.de)

21 **Abstract.** Earth System Models (ESMs) are a key tool for studying the climate under changing conditions. Over recent  
22 decades, it has been established to not only rely on projections of a single model but to combine various ESMs in multi-model  
23 ensembles (MMEs) to improve robustness and quantify the uncertainty of the projections. The data access for MME studies  
24 has been fundamentally facilitated by the World Climate Research Programme's Coupled Model Intercomparison Project  
25 (CMIP) - a collaborative effort bringing together ESMs from modelling communities all over the world. Despite the CMIP

26 standardization processes, addressing specific research questions using MMEs requires unique ensemble design, analysis, and  
27 interpretation choices. Based on the collective expertise within the Fresh Eyes on CMIP initiative, mainly composed of early-  
28 career researchers engaged in CMIP, we have identified common issues and questions encountered while working with climate  
29 MMEs. Here, we provide a comprehensive literature review addressing these questions. We provide statistics tracing the  
30 development of the climate MMEs analysis field throughout the last decades, and, synthesizing existing studies, we outline  
31 guidelines regarding model evaluation, model dependence, weighting methods, and uncertainty treatment. We summarize a  
32 collection of useful resources for MME studies, we review common questions and strategies, and finally, we outline emerging  
33 scientific trends, such as the integration of machine learning (ML) techniques, single model initial-condition large ensembles  
34 (SMILEs), and computational resource considerations. We thereby aim to support researchers working with climate MMEs,  
35 particularly in the upcoming 7th phase of CMIP.

## 36 **1 Introduction**

37 Earth system models (ESMs) are a key tool for assessing the future climate under changing conditions. Starting from the  
38 seminal work of Manabe and Hasselmann (e.g. Manabe and Strickler, 1964; Manabe and Bryan, 1969; Manabe and Wetherald,  
39 1967; Hasselmann, 1976), who were awarded the 2021 Nobel Prize in Physics, climate models have continuously evolved  
40 over decades. During this process, models have become progressively more complex, encapsulating processes related to  
41 aerosols, atmospheric chemistry, the carbon cycle, and ocean biogeochemistry. This evolution occurred in parallel with  
42 advances in Earth system observations, high-resolution numerical models giving insight into smaller-scale phenomena (e.g.,  
43 detailed radiative transfer models, cloud-resolving models, large-eddy simulations), and growing computational power (e.g.  
44 Gettelman et al., 2022; Schneider et al., 2017) allowing model resolution to steadily improve.

45 Since the beginning of large-scale atmospheric modelling, intercomparisons among models have been carried out. Initially,  
46 this intercomparison was mostly performed for numerical weather prediction as computational resources limited the  
47 intercomparison of studies in the climate context, and a clear experimental strategy was lacking (Gates, 1992). Since the 1970s,  
48 the Working Group on Numerical Experimentation (WGNE), supporting the World Climate Research Programme, has  
49 organized several intercomparison projects among climate models. The first international systematic intercomparison  
50 framework for climate models was established in 1990 in the context of the Atmospheric Model Intercomparison Project  
51 (AMIP; Gates, 1992). In the early 1990s, the Intergovernmental Panel on Climate Change (IPCC) provided an intercomparison  
52 of atmospheric models in their first assessment report (AR; Gates, 1992). In the following years, Räisänen (1997) advocated  
53 the need for quantitative model comparison and raised the thought that the agreement between models can indirectly serve as  
54 a measure for the reliability of the simulations. Accordingly, Räisänen and Palmer (2001) introduced a probabilistic perspective  
55 on MME projections. The authors quantified the probability of specific climate events happening based on 17 coupled  
56 atmosphere-ocean general circulation models (AOGCMs). Contemporaneously, AMIP was followed by the Coupled Model

57 Intercomparison Project (CMIP) coordinated by the World Climate Research Programme (WCRP), which also incorporated  
58 results from AOGCMs (Meehl et al., 2000).

59 While the first phase of CMIP was limited to control runs, new standardized scenarios were incorporated throughout the phases  
60 of CMIP with an increasing number of international model centres contributing simulations. Concurrently, the volume of data  
61 has been steadily increasing (Williams et al., 2016) and is stored within a standardized format at the Earth System Grid  
62 Federation (ESGF) central repository (Cinquini et al., 2012). In more recent CMIP generations, a variety of supporting  
63 experiments is conducted (e.g. Eyring et al., 2016), including paleoclimate runs (simulations of the ‘distant past’), historical  
64 runs (simulations of the ‘recent past’), control runs to study natural variability, as well as various experiments. Finally, future  
65 climate change simulations are performed for various greenhouse gas emission scenarios such as abrupt carbon dioxide  
66 doubling or quadrupling to derive climate sensitivity (measure of how much the Earth's climate system will warm under a  
67 doubling of atmospheric CO<sub>2</sub> concentration), as well as for multiple plausible future emission scenarios (O’Neill et al., 2017;  
68 Riahi et al., 2017; van Vuuren et al., 2025). The latter denote diverse scenarios of evolution of the global society (including  
69 population, economy, and technology) which thus lead to differing emissions of greenhouse gases (CO<sub>2</sub>, CH<sub>4</sub>, NO<sub>2</sub>) and other  
70 air pollutants until the end of the 21st century and are associated with different climate change mitigation and adaptation  
71 policies and challenges (IPCC, AR6). These CMIP projections have proven essential for informing mitigation and adaptation  
72 strategies to climate change at the global and regional scales (Meehl et al., 2000). For regional analysis, the CMIP output is  
73 often downscaled to finer resolution, e.g. by using the CMIP output as boundary conditions for regional climate models  
74 (RCMs). This is done e.g. in the WCRP COordinated Regional climate Downscaling EXperiment (CORDEX), which provides  
75 a coordinated framework for producing and evaluating regional climate projections across multiple domains worldwide  
76 (Giorgi, 2019; Gutowski Jr. et al., 2016)

77 The main components of an ESM describe the atmosphere, ocean, cryosphere, land, and increasingly, the carbon cycle and  
78 other biogeochemical processes. Each component involves a variety of interacting phenomena occurring at a wide range of  
79 spatial and temporal scales (e.g. Gettelman et al., 2022). In all ESMs, the continuous behavior of the atmosphere and ocean is  
80 first discretized in space and time via the so-called “model dynamical core” which encompasses known governing equations  
81 that capture resolved (grid-scale) phenomena and parameterization schemes that represent unresolved or poorly understood  
82 (subgrid-scale) processes. ESMs differ in the choice of computational grids (e.g., latitude-longitude structured grids,  
83 icosahedral grids, variable resolution cube-sphere grids), numerical methods for solving the dynamical core equations, and in  
84 parameterization schemes. While some parameterizations are based on well-established physical theory, others, particularly  
85 those related to clouds or turbulence, remain subject to substantial uncertainty. In addition, computational limitations restrict  
86 the accuracy with which models can represent certain relevant processes. Therefore, the decisions made at modelling centers  
87 make each ESM an imperfect attempt to represent a multitude of highly complex, nonlinear processes, and the synchronized

88 interplay among them. Depending on the interest of the end user, some of these necessary idealization decisions may be more  
89 suitable than others.

90 Combining several ESMs to multi-model ensembles (MMEs) can have numerous advantages compared to individual  
91 simulations, e.g. to account for the uncertainty arising from the differing modelling decisions (model uncertainty). Starting in  
92 the weather forecasting community, numerous studies have shown the benefits of ensemble predictions compared to  
93 predictions based on single models (Doblas-Reyes et al., 2003; Krishnamurti et al., 1999), e.g. the North American MME  
94 showed improvements in various skill metrics (correlation, RMSE, RPSS, and reliability) compared to individual models used  
95 before (Kirtman et al., 2014). Inspired by these findings, studies in the climate context also analyzed the potential benefits  
96 from working with MMEs for projections. In climate model evaluation, the MME projections have proven to outperform  
97 individual model projections in numerous studies e.g. regarding the mean (Gleckler et al., 2008; Knutti et al., 2010a; Lambert  
98 and Boer, 2001; Palmer et al., 2005; Phillips and Gleckler, 2006; Pincus et al., 2008; Reichler and Kim, 2008) and variability  
99 (Zhang et al., 2007). The enhancement of the signal and cancellation of errors contribute to these advantages (Doblas-Reyes  
100 et al., 2005; Hagedorn et al., 2005; Smith et al., 2013). Becker et al. (2022) highlight the practical advantage of the continuous  
101 operation of MMEs, which can be maintained even when individual modelling centers are temporarily unable to contribute,  
102 for example due to technical or political constraints. They further provide an example where the use of a MME enabled the  
103 identification of outlier behavior in ENSO predictions, which could subsequently be traced back to previously unknown  
104 deficiencies in the underlying reanalysis dataset, thereby supporting the model improvement. Furthermore, an ensemble  
105 approach reduces the risk of selecting a model outlier with particularly large biases.

106 Given these benefits, MME projections have become an established tool for climate studies addressing a broad range of  
107 research questions, also being the standard method to analyze and present results in the Assessment Reports (ARs) of the  
108 Intergovernmental Panel on Climate Change (IPCC), where the state-of-the-art knowledge on climate change is reviewed. For  
109 researchers, MMEs provide an efficient way to get an overview of general tendencies for specific questions. Also for non-  
110 experts, presenting results in a synthesized format as e.g. in the context of MME also facilitates accessibility and interpretation  
111 (Knutti et al., 2010a), underlining the benefits of MMEs for the users.

112 It is important to recognize that CMIP constitutes an “ensemble of opportunity” (Tebaldi and Knutti, 2007; Sanderson et al.,  
113 2012; Merrifield et al., 2023), as it reflects the collection of readily available simulations rather than a systematically designed  
114 sample. Contributing institutions range from long-established, well-resourced climate modelling centers to newer groups with  
115 sufficient computational resources to run adapted versions of existing models. While this inclusivity broadens participation,  
116 such ensembles of opportunity are not designed to constitute a statistically representative sample of multi-model uncertainty  
117 (Merrifield et al., 2023). In this context, the superiority of MMEs is not universal. There are cases in which individual models  
118 can outperform the ensemble mean, for instance when the averaging inherent to MMEs suppresses relevant signals that are

119 well represented in only a subset of models. This can occur for specific physical processes, or extremes, where ensemble  
120 averaging may smooth physically meaningful variability or dampen circulation-driven responses. Moreover, if most models  
121 in an MME share common structural components, parameterizations, or tuning strategies, systematic biases can persist in the  
122 ensemble mean. In such cases, individual models with alternative formulations may provide more accurate representations for  
123 specific variables, regions, or applications.

124 The availability of standardized climate model outputs facilitated model intercomparison and has naturally inspired the use of  
125 MMEs since the beginning of the 2000s (Tebaldi and Knutti, 2007). Consequently, the AR3 of the IPCC (2001) presented  
126 many results based on MME means, accompanied by measures of inter-model variability (Tebaldi and Knutti, 2007). In the  
127 AR4 of IPCC (2007), model projections were only included if the models were successors from previous generations, thus a  
128 model selection *de facto* has taken place (Knutti et al., 2010b). To support IPCC lead authors for the AR5 and later, a “Good  
129 Practice Guidance Paper” was published in 2010, summarising current recommendations for the work with MMEs (Knutti et  
130 al., 2010b).

131 In the meantime, numerous studies have proposed diverse methods for MME studies. However, it is challenging to have an  
132 overview of these studies, and there is still a lack of guidelines on how to combine models within MMEs (Herger et al., 2018).  
133 The design of MME studies involves a set of decisions related to model selection, weighting, and uncertainty measures. Each  
134 of these decisions requires careful consideration of a broad range of aspects and often entails compromises that differ  
135 depending on the research question. This individuality makes it challenging or even impossible to establish universally  
136 applicable guidelines for MME studies. However, we believe it is valuable to give an overview of the key aspects to consider,  
137 and in some cases, present approaches that the Fresh Eyes on CMIP community has found to be useful. With this, we hope to  
138 support researchers that have newly entered the field of climate science, but also to provide an overview of existing resources  
139 and approaches for more experienced scientists, particularly for (but not restricted to) the upcoming 7th phase of CMIP.

140 While the focus of this paper is on the challenges associated with combining various ESMs within a MME, it should be pointed  
141 out there are other types of climate ensembles. Besides such uninitialized simulations, there are initialized climate model  
142 ensembles that are routinely used for seasonal prediction (see e.g. Becker et al., 2020, 2022; Buontempo et al., 2022; Kirtman  
143 et al., 2014; Min et al., 2025). Initialized climate model ensembles are based on accurate initialization and thus have an  
144 emphasis on assimilation procedures to capture the atmosphere, ocean and land conditions. While their goals differ from those  
145 of CMIP, initialized prediction ensembles face similar challenges related to ensemble design, model weighting, and evaluation  
146 against observations. Further ensemble types include initial condition ensembles (ICEs) and perturbed parameter ensembles  
147 (PPEs) (IPCC, AR5). ICEs are generated with a single climate model using varying initial conditions (i.e., perturbed initial  
148 state) to address the uncertainty due to natural or internal variability. If sufficiently many ensemble members are available,  
149 they are referred to as Single Model Initial-condition Large Ensembles (SMILEs). The perturbed parameter ensemble (PPEs)

150 also compares multiple realizations from a single climate model, but in this case, a set of chosen physical parameters which  
151 are assumed to affect the quantity of interest (e.g., global mean surface temperature) is systematically varied to quantify the  
152 effect on model outcome (e.g. Eidhammer et al., 2024; Sexton et al., 2021). This enables a systematic exploration of intra-  
153 model uncertainty. Finally, the so-called grand ensembles are based on a combination of various ensemble types (IPCC, AR6).

154 In the following section, we conduct a comprehensive literature review on studies regarding model evaluation (2.1), model  
155 dependence (2.2), model selection and weighting methods (2.3) and uncertainty characterization (2.4). In this context, we also  
156 provide a summary of useful tools for MME analysis (2.5). In Section 3, we complement these guidelines with a collection of  
157 frequently occurring topics and challenges based on the experience of the WCRP Fresh Eyes on CMIP community. In Section  
158 4, we discuss emerging trends for working with MMEs such as machine learning (ML), SMILEs and the necessity for more  
159 awareness of computational resources associated with MME studies.

## 160 **2 Guidelines for working with MMEs**

161 Over 84 General Circulation Models (GCMs) from at least 43 international institutes are available through CMIP ([https://wcrp-](https://wcrp-cmip.org/map/)  
162 [cmip.org/map/](https://wcrp-cmip.org/map/)). When addressing any research question, the need for specific variables, scenarios, resolutions or experiments  
163 narrows the pool of available models. However, the remaining number is often still large, prompting the following questions:  
164 Should all available models be used, or only a subset? How can the models be identified that are most suitable for such a  
165 subset? The two primary objectives when selecting models are to optimize model performance and to reduce duplicated  
166 information (Herger et al., 2018). As adequate selection criteria are central to the design of MME studies, we aim to provide  
167 guidance for the choice of models in this section.

### 168 **2.1 Model Evaluation**

169 Model evaluation refers to the systematic assessment of climate model simulations against observational reference data in  
170 order to compare model performance and identify biases. For an overview of model bias see Appendix B. In practice, this  
171 involves benchmarking historical simulations with respect to observed climate statistics, such as mean states, variability, spatial  
172 patterns, and relevant physical processes.

### 173 **Observation Datasets for Model Evaluation**

174 Observational reference datasets used for model evaluation include both direct observations and reanalysis products.  
175 Reanalysis datasets are physically consistent products produced by assimilating diverse observational data into a numerical  
176 weather or climate model. They combine the broad spatial and temporal coverage of models with observational constraints  
177 and are therefore widely used as reference datasets. Direct observations include paleoclimate data, ground-based measurements  
178 over land and ocean (e.g., ships, buoys and sail drones), aircraft and balloon measurements, and satellite data. Paleoclimate

179 data give insight into the state of the Earth’s climate hundreds to millions of years ago, offering valuable constraints for  
180 paleoclimate simulations that help us understand recent and future climate change in the context of longer-term climate  
181 variability. For the more recent past, most of the reference observations originated in land in-situ measurements, which are not  
182 equally distributed around the globe (e.g., there are more land measurement stations in the Northern Hemisphere than in the  
183 Southern Hemisphere). The advent of Earth observation satellites has revolutionized the availability and coverage of global  
184 reference datasets. However, satellite datasets are limited to the time after the 1970s, depending on the variable of interest.

185 All these datasets have distinct advantages and disadvantages: They encompass different spatial and temporal scales, cover  
186 different locations and time periods, rely on different measurement techniques, or vary in accuracy. See e.g. Sippel et al.,  
187 (2024) for challenges in observational data. Associated uncertainties also differ, e.g. due to instrument uncertainty, calibration  
188 limitations, or interpolation procedures. Accounting for these uncertainties in the reference datasets can be done by combining  
189 multiple datasets (Notz et al., 2016). It also facilitates signal detection for subsequent comparison with model ensemble outputs  
190 (Santer et al., 2008). Observational ensembles have been paired with MMEs in studies e.g. with regard to the tropical  
191 troposphere (Santer et al., 2008) or to Antarctic sea ice (Roach et al., 2018). Depending on the variable of interest, commonly  
192 used reanalysis datasets are ERA5 (produced by ECMWF), MERRA-2 (produced by NASA GSFC), NCEP DOE R-2  
193 (produced by NOAA), JRA-3Q (produced by Japan Meteorological Agency).

194 Moreover, model evaluation using observations is not always straightforward, as observational sensors do not necessarily  
195 measure variables simulated by climate models. To ensure an “apple-to-apple comparison”, observed quantities must be  
196 converted into model-output-like variables, or vice versa. For example, software has been developed which enables simulating  
197 what a satellite would observe over the model atmosphere. Moreover, it must be assured that observations and simulations  
198 have the same temporal and spatial resolution, including the horizontal grid and number of vertical levels (Simpson et al.,  
199 2025), which can be achieved by appropriate regridding methods. See Section 3.5 for details on regridding.

200 Generally, there are two approaches to model evaluation: (i) The performance-oriented approach focuses on identifying the  
201 models whose output is closest to observations or reanalysis data. (ii) The process-oriented approach seeks models that best  
202 capture the dynamics of interest. Regardless of the approach, it is essential for any research project to report on the performance  
203 of all models available before applying any ranking or weighting methods, and the selection criteria should be reported  
204 transparently (Knutti et al., 2010a). Such evaluations are sometimes already available in the literature and can be referenced.  
205 But in that case it is important to make sure that they cover the variables, scales, and other factors relevant to the specific  
206 research questions.

## 207 **Performance-oriented Evaluation**

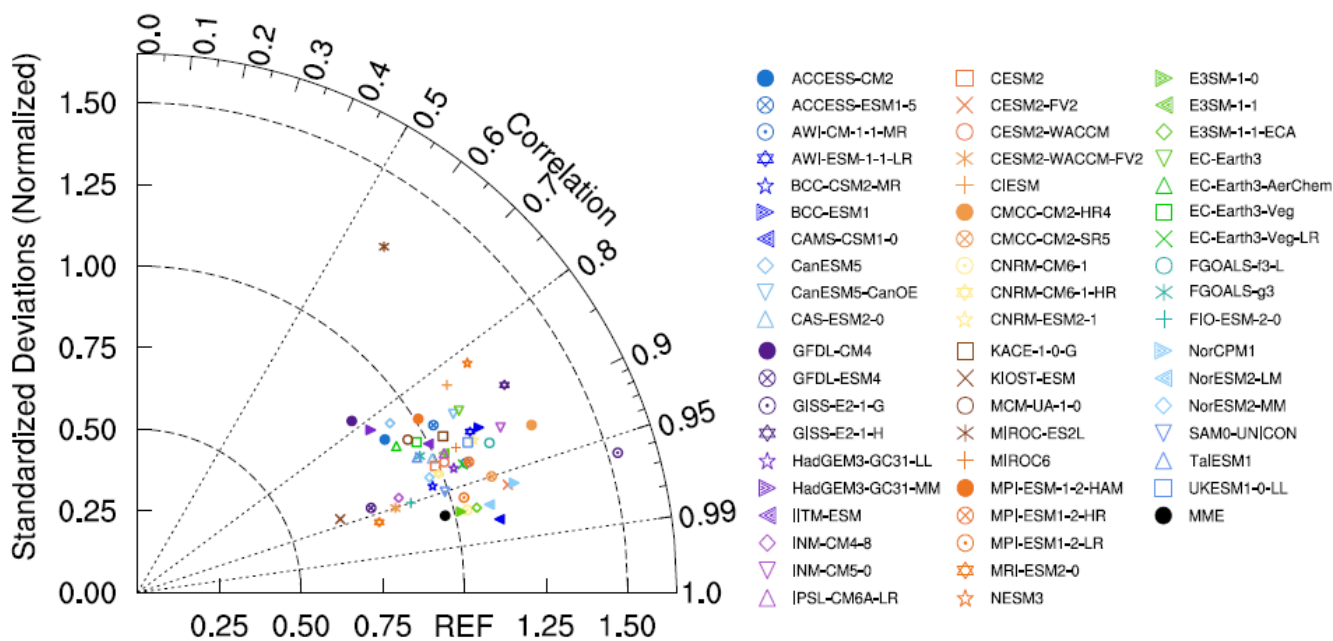
208 For shorter timescale forecasts, predictions can be verified within days as observations become available. This is typical of  
209 weather forecasting and initialized climate model simulations, in which models are started from observation-constrained initial  
210 conditions. Such near-term verifiability offers an opportunity to build confidence in models, particularly for climate services  
211 and decision-relevant applications. Although initialized and uninitialized climate projections address different time horizons,  
212 linking insights from both may help contextualize uncertainties and enhance trust in long-term projections. Climate projections  
213 addressing longer time scales cannot be directly verified in real time, as the relevant time scales (decades to centuries) preclude  
214 immediate verification. This is the case for uninitialized climate model simulations, which represent the standard approach for  
215 long-term climate projections and are the focus in this review. Accordingly, climate model performances are evaluated with  
216 reference to past and present-day climatology (Knutti, 2010).

217 Performance-oriented model evaluation is based on the assumption that models that fail to perform well for the past regarding  
218 some specific climate phenomena will also do so for the future. While this assumption is commonly accepted, it also is a  
219 limitation of this approach as the role of specific circulation patterns and their interactions might change throughout the 21st  
220 century. In this context, Knutti et al., 2010a found that model performance evaluated for the past correlates only weakly with  
221 the magnitude of the projected change in the future, illustrating that constraining models based on past performance does not  
222 necessarily reduce future inter-model spread. Given these pitfalls, Mendlik and Gobiet (2016) propose to only remove the  
223 severely unrealistic models. A detailed assessment on how to deal with outliers can be found in Subsection 3.3.

224 Because uninitialized climate model simulations are free-running and not constrained by observations, performance-oriented  
225 evaluation cannot rely on a direct comparison of individual events or temporal trajectories. Instead, model evaluation is  
226 necessarily based on climatological characteristics, such as mean states or spatial patterns. Evaluating this climatological  
227 performance comes down to the choice of appropriate metrics. Model ranking has been found to be sensitive to this choice  
228 (Gleckler et al., 2008). However, for specific variables, the model projections may be largely independent of the choice of  
229 underlying metrics and ranking methods (Santer et al., 2009). Given the diversity of possible research questions, there is no  
230 single or combined performance metric that can reliably identify the “best” model independent of the research question. While  
231 this may sound disappointing since it prevents the standardization of model evaluation, it also has the advantage of reducing  
232 the effect of model convergence due to tuning (Knutti, 2010), which allows for a more reliable representation of future  
233 uncertainty and decreases the likelihood of making overconfident predictions. Generally, a metric is recommended if it’s as  
234 simple as possible while at the same time being as statistically robust as possible, meaning that the dependence on  
235 specifications of the metric is rather low (Knutti et al., 2010b). Therefore, for any study, it is essential to use metrics that are  
236 relevant to the specific research question while also matching the spatial and temporal scale of the phenomenon in question.

237 Taylor diagrams (Taylor, 2001) have become a widely used tool to visualize performance-oriented model evaluation, helping  
238 to identify better performing models as well as outliers. They are applied across the full range of climate-related topics,

239 including e.g. the Indian Summer Monsoon (Roy et al., 2019) and seasonal mean temperatures (Tang et al., 2016). In a Taylor  
 240 diagram, the radial distance from the origin represents the model standard deviation, the angle from the horizontal axis encodes  
 241 the correlation with observations, and the geometric distance to the reference point (defined as the observed standard deviation  
 242 and correlation = 1) equals the centered root mean square error, quantifying pattern mismatch after mean removal. Models  
 243 closer to the observed standard deviation, along with higher correlation coefficients (and therefore lower root mean square  
 244 error) are considered as better performing models for specific climate features (Taylor, 2001). The angular/azimuthal position  
 245 in Taylor diagrams represents the pattern correlation coefficient between CMIP6 models and observations, while the radial  
 246 distance indicates the ratio of the standard deviation of CMIP6 models to that from an observational data set. For example, the  
 247 Western Pacific pattern, a prominent teleconnection pattern during the boreal winter over the North Pacific, was analyzed for  
 248 56 CMIP6 models using a Taylor diagram (Fig. 1, Aru et al., 2023). It depicts that the spatial correlations of the geopotential  
 249 height anomalies at 500-hPa over the Western North Pacific between individual CMIP6 models and observations exceed 0.6.  
 250 In reproducing spatial patterns, the mean of the MME outperforms most individual models, which is evidenced by a spatial  
 251 root mean square deviation of 0.97. This diagram also makes it possible to identify outlier models, such as the MIROC-ES2L  
 252 in this example. Selecting only the best performing models can improve the final MME mean.



253

254 Fig. 1. Example for the use of a Taylor diagram showing the geopotential height anomalies  
 255 at 500-hPa over the Western North Pacific (20°N–80°N, 120°E–120°W) in individual CMIP6  
 256 models, MME and observations, taken from Aru et al. (2023).

257 A frequent challenge in climate model evaluation is determining whether models yield correct results for incorrect reasons,  
258 due to compensating errors (Eyring et al., 2016; Ivanova et al., 2016). There is a possibility that, while a model appears to  
259 accurately represent some variable, the underlying processes are not well-captured, which could mask inherent biases in the  
260 model. For example, analysing CMIP6 models, Zhao et al. (2022) reported that the cloud radiative effect reveals compensating  
261 errors between the modeled total cloud fraction and the liquid water path. These errors offset each other, resulting in a smaller  
262 net error in the cloud radiative effect. Di Luca et al. (2020a) addressed the issue of error compensation in CMIP5 simulations  
263 of hot temperature extremes by developing a new error metric called the “additive error.” This metric adds up the absolute  
264 errors of four components contributing to temperature extremes: the long-term mean, seasonality, diurnal temperature range,  
265 and the local temperature anomaly on the day of the extreme. Compared to traditional bias or absolute error metrics, the  
266 additive error more sensitively captures the total error in extreme temperature estimates. Furthermore, Di Luca et al. (2020b)  
267 defined a new error estimator that aims to minimise error compensation.

268 It is important to remember that models are calibrated with the aim to reduce anomalies compared to observational data before  
269 becoming available in new CMIP generations. During this calibration (often referred to as tuning), parameters, typically  
270 associated with unresolved processes such as clouds, convection, or boundary-layer dynamics, are adjusted to improve  
271 agreement with observations. Consequently, improvements in overall model performance in new CMIP generations do not  
272 necessarily stem from enhanced capabilities in capturing relevant processes, but may instead result from optimized calibration  
273 (Knutti, 2010). A related issue is that the same observational datasets used for model calibration are often also employed for  
274 model evaluation, which is not optimal as calibration and evaluation datasets ideally should be independent. This concern is  
275 even more pronounced when using reanalysis products as reference data, since climate models are an integral part of their  
276 generation.

277 Additionally, observational data can influence model performance through the forcings themselves. For example, in  
278 concentration-driven CO<sub>2</sub> simulations, observed atmospheric concentrations are prescribed directly for historical simulations,  
279 rather than being computed from emissions, as in emission-driven models. This approach further constrains the model output,  
280 since the model does not simulate atmospheric CO<sub>2</sub> concentrations from emissions via an interactive carbon cycle.  
281 Consequently, apparent improvements visible in the model’s evaluation do not necessarily indicate a better representation of  
282 the carbon cycle itself.

283 Ideally, the evaluation process also allows insights on how well basic dynamic processes relevant to the research questions are  
284 reproduced in models (Knutti et al., 2010b). For a research question regarding rainfall, for example, this could mean to not  
285 only analyze the precipitation pattern, but also inspect wind patterns to see if the associated circulation is captured well.  
286 Process-oriented model evaluation specifically targets the model performance concerning such dynamics.

## 287 **Process-oriented Evaluation**

288 This evaluation approach shifts from traditional performance-oriented evaluation to more detailed, process-oriented metrics,  
289 which are critical for advancing the next generation of ESMs. Eyring et al. (2005) and Gleckler et al. (2008) emphasise the  
290 need to evaluate a wide range of climate processes, since accurately simulating one aspect does not ensure accuracy in others.  
291 These authors initiated the development of a comprehensive set of model metrics to assess important processes in climate  
292 simulations. Process-oriented evaluation identifies sources and limitations of predictability, guiding model development by  
293 revealing deficiencies in the representation of physical processes and thereby enhancing the reliability of climate projections  
294 (Eyring et al., 2016). By incorporating process-oriented analysis into diagnostic packages (examples in Subsection 2.5),  
295 evaluations become reproducible, accelerating model improvements and establishing benchmarks for progress. As with any  
296 standardization effort, however, such benchmarks must be applied with care, as they have the potential to promote model  
297 similarity. Another relevant resource in the context of process-oriented evaluation is Simpson et al. (2025), who review the  
298 ability of climate models to reproduce historically observed forced trends and outline best practices for confronting modeled  
299 and observed signals. Within the MME framework, process-based approaches help identify which processes contribute most  
300 to inter-model differences and provide insights into the mechanisms behind model performance. Here, we outline some  
301 common use cases and techniques.

302 *Process-oriented diagnostics to reduce model bias:* One major focus in the development of process-oriented metrics is the  
303 investigation of phenomena with strong bias in the models, as e.g. the Madden-Julian Oscillation (MJO), the dominant mode  
304 of tropical intraseasonal variability. To better understand the origins of these biases, a number of diagnostics has been  
305 developed to facilitate improvements in the representation of the MJO in weather and climate models (Ahn et al., 2020; Li et  
306 al., 2022; Wang et al., 2020). The first process-oriented multi-model comparison study on MJO teleconnections found that  
307 biases in simulating the position of the Pacific westerly jets, together with deficiencies in MJO representation, contribute  
308 substantially to errors in MJO teleconnections (Ahn et al. 2017, Henderson et al. 2017). Similar efforts exist for the El Niño–  
309 Southern Oscillation (ENSO), for which Planton et al. (2021) provide a dedicated metrics package.

310 *Improving projections by process-oriented multiple diagnostic ensemble regression:* Karpechko et al. (2013) developed the  
311 multiple diagnostic ensemble regression (MDER) method that constrains climate projections using observed diagnostics,  
312 applying it to Antarctic ozone columns. By identifying key processes that influence ozone, MDER explains a substantial  
313 fraction of the inter-model spread in projected ozone across climate chemistry models and outperforms the unweighted multi-  
314 model mean in pseudo-realistic validation. Building on this approach, Wenzel et al. (2016) applied the MDER algorithm,  
315 implemented as a diagnostic in ESMValTool (see Subsection 2.5), to analyze projections of the austral jet position under the  
316 RCP4.5 scenario in CMIP5 simulations. They found that MDER reduces uncertainty in the ensemble-mean projection without  
317 substantially altering the long-term mean position of the jet.

318 *Identifying the role of model configurations:* Another significant aspect of process-oriented model evaluation is understanding  
319 how specific characteristics are influenced by model configurations, such as resolution and parameterization schemes. Kim et  
320 al. (2018) proposed a set of diagnostics to assess how model physics affect the representation of tropical cyclones, particularly  
321 their intensity in GCMs. The findings suggest that model-specific factors, beyond large-scale environmental parameters, play  
322 a key role in shaping tropical cyclones' intensity, with differences in convection schemes contributing significantly to the inter-  
323 model spread. Wing et al. (2019) and Moon et al. (2020) further applied these methods, with Moon et al. (2020) showing that  
324 tropical cyclone wind structures are strongly influenced by model resolution. Accordingly, Dirkes et al. (2023) emphasizes the  
325 necessity of applying the developed diagnostics for tropical cyclone analysis in CMIP6 models.

326 *Using idealization or a hierarchy of models:* Another approach is to design model configurations that isolate individual  
327 processes and components, allowing to test their relevance for specific phenomena. For example, Katzenberger et al. (2024)  
328 employed an aquaplanet configuration with a circumglobal land stripe to evaluate the meridional circulation, particularly the  
329 Hadley cell, in an idealized setup. By shifting the landstripe north and southwards, and by modifying the surface albedo or  
330 aerosol concentrations, the role of these features in shaping monsoon dynamics could be systematically isolated. More  
331 generally, iteratively adding components and increasing the complexity and realism of the setup within a hierarchy of models  
332 enables the isolation of individual processes and the assessment of their contributions to the overall model performance. See  
333 also e.g. Zhou and Xie (2018) for more insights to this approach.

334 *Using causal inference:* In Section 4.1, we provide insights into how ML techniques can be applied to improve process-based  
335 evaluation by identifying causal relationships.

336 Another example of process-oriented assessment is provided by Fasullo et al. (2020), who present a thorough analysis of CMIP  
337 representation of the leading Earth system modes of variability. Additional applications include regime-based evaluation  
338 approaches of low-level marine clouds, where distinguishing stratocumulus from shallow cumulus regimes has helped  
339 diagnose persistent cloud-cover and radiative biases in CMIP6 and CMIP5 models and inform targeted model improvements  
340 (Črnivec et al., 2023; Cesana et al., 2023). Process-based analyses have also demonstrated that the ENSO–Indian Summer  
341 Monsoon teleconnection is robustly represented in CMIP5 and CMIP6 models, consistent with a realistic simulation of the  
342 coupled Hadley–Walker circulation and associated precipitation responses (Roy and Tedeschi, 2016; Roy et al., 2017; Fasullo  
343 et al., 2020). We provide further details and examples of process-oriented analyses in the Appendix C.

## 344 **2.2 Model Dependence**

345 ESMs are developed by multiple modelling groups worldwide. Ideally, the models in a MME would be independent, thereby  
346 providing an adequate representation of the epistemic uncertainty. Historically, climate projections are derived by calculating  
347 simple averages across the MME, based on the assumption that the ensemble mean offers the most accurate representation of

348 the Earth system by synthesizing the collective modelling efforts (Abramowitz et al., 2019; Knutti et al., 2010a). Assuming  
349 independence implies that the MME reflects a sufficiently broad range of uncertainties, and the averaging smooths out  
350 individual model biases. In practice, however, the development of ESMs is often not independent (Pincus et al., 2008).

351 Components that address modelling challenges or have demonstrated strong performance are often shared among multiple  
352 ESMs, including e.g. the dynamical core for resolving grid-scale dynamics or components addressing sub-grid-scale  
353 phenomena (e.g., parameterization schemes). For example, the McICA radiation scheme (Pincus et al., 2003) provides an  
354 efficient and flexible representation of one-dimensional radiative transfer in a cloudy atmosphere, and is thus implemented in  
355 multiple ESMs such as several US models (NSF NCAR CESM2, NOAA GFDL-CM4, DOE E3SM-1-0), the Canadian model  
356 (CanESM5), the UK model (HadGEM3), and the Norwegian model (NorESM2). Similarly, the NEMO ocean model is widely  
357 used across modelling centers, including e.g. HadGEM3 and NorESM2, further underscoring the sharing of model  
358 components. Fig. 2 illustrates the shared model history tracing back to a few AGCMs (Kuma et al., 2023).

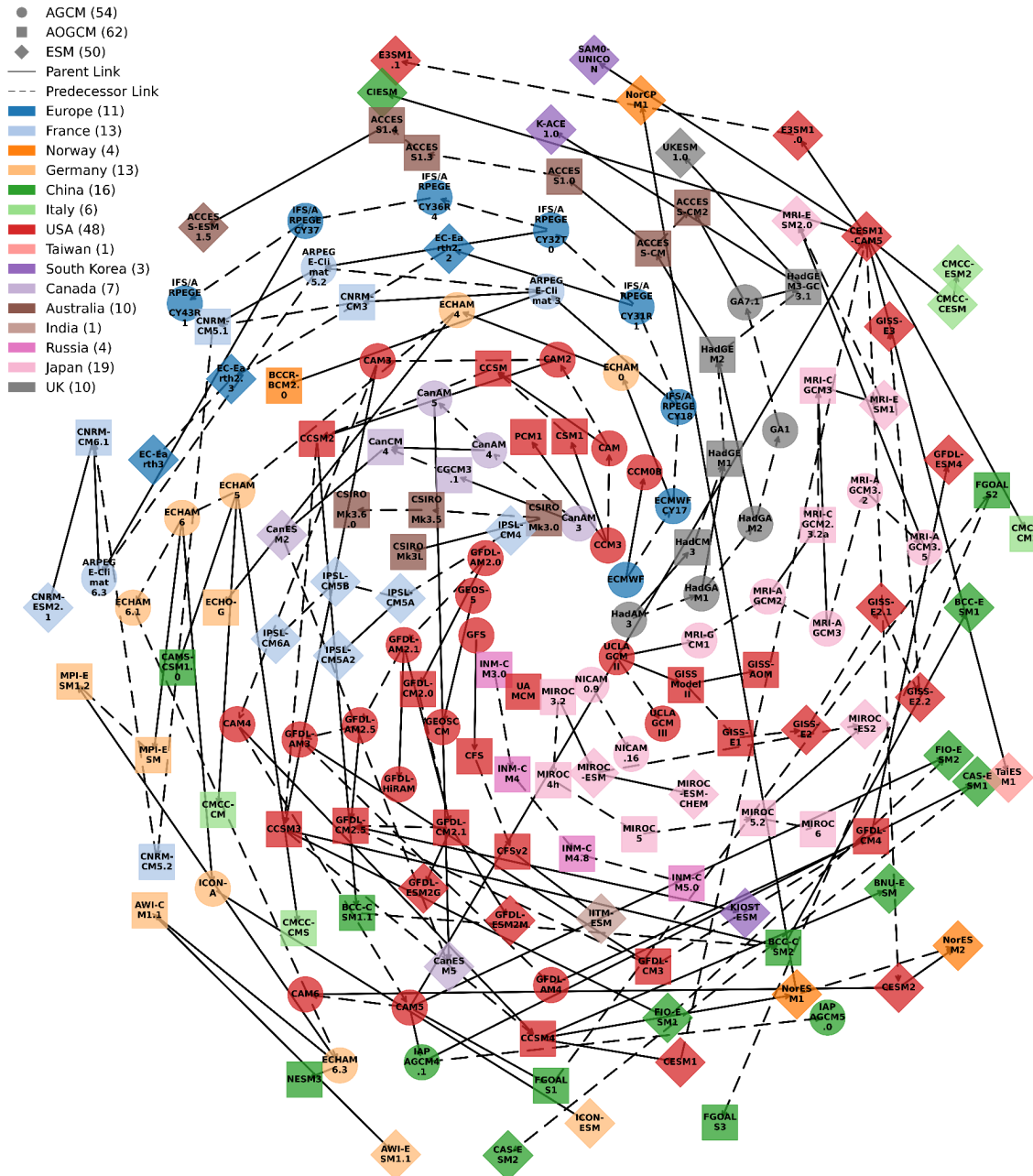
359 In addition to dependencies between modelling groups, individual centers often contribute multiple closely related model  
360 configurations, for example differing in horizontal resolution (e.g., MPI-ESM1.2-HR for high resolution versus MPI-ESM1.2-  
361 LR for low resolution) or in the inclusion of additional components, such as interactive vegetation in EC-Earth3-Veg compared  
362 to EC-Earth3. If such dependencies are not accounted for and all models are included with equal weight in a multi-model  
363 mean, modelling centers that provide several related configurations effectively receive greater weight than others. This issue  
364 is particularly relevant in multi-model studies with a limited number of included models.

365 Diverse analyses confirm that the number of independent models in CMIP is smaller than the total number of participating  
366 models (Jun et al., 2008; Masson and Knutti, 2011; Pennell and Reichler, 2011). Because errors across different models are  
367 often being correlated (Knutti et al., 2010a), the lack of independence can lead to amplified biases (Jun et al., 2008; Knutti,  
368 2008; Reichler and Kim, 2008; Tebaldi and Knutti, 2007). Moreover, apparent convergence among model results and the  
369 associated reduction in ensemble uncertainty may be mistakenly interpreted as strong agreement between models, when in fact  
370 they arise from structural dependencies.

371 The lack of a universally accepted and unambiguous definition of model independence complicates efforts to systematically  
372 account for model dependence in MME studies. Some definitions focus on the conceptual idea of whether or not a model adds  
373 novel additional information to the MME (Masson and Knutti, 2011). Others adopt a more analytical approach to  
374 understanding model dependence, offering examples for evaluating model dependence and using their framework (e.g., Annan  
375 and Hargreaves, 2017). Despite such advances, no broadly accepted solution has yet emerged. Further approaches, such as  
376 weighting schemes (Subsection 2.3) have been proposed, but these tend to be problem-specific and struggle to capture the full  
377 complexity of model dependencies. The metadata reporting requirements introduced in CMIP6 have made comprehensive  
378 assessments of model dependence possible, thereby representing a meaningful advance in transparency. As new model

379 generations are developed and incorporated to CMIP, continued efforts to quantify and correct for model dependence will be  
 380 essential to ensure robust ensemble projections that reflect true uncertainty.

381



382

383 Fig. 2. Spiral plot of climate model dependencies, adapted from Kuma et al. (2023). The oldest model in any given family is  
384 in the center of the plot, spiralling out as more models are made. Model type is differentiated by shape of marker, and link type  
385 is differentiated by arrow type (solid for parent or dashed for predecessor). Models developed in different countries are assigned  
386 distinct colors. Markers indicate atmosphere general circulation models (AGCMs), atmosphere-ocean global circulation  
387 models (AOGCMs), and Earth system models (ESMs). Numbers of models from each country are indicated in brackets in the  
388 legend. ECMWF models are denoted by the country “Europe”.

### 389 **2.3 Model Selection and Weighting Methods**

390 CMIP MME weighting and selection techniques are used to categorize the CMIP models based on historical model  
391 performance (see Subsection 2.1) and independence (see Subsection 2.2) using several metrics (Palmer et al., 2023), which is  
392 crucial for optimizing accuracy and reliability in projections (Strobach and Bel, 2020). Several performance-based and  
393 statistical approaches are used for MME weighting (Bhowmik and Sankarasubramanian, 2020; Brunner et al., 2020).  
394 Performance-based weighting assigns weights based on the ability to reproduce observed historical climate patterns, while  
395 statistical model weighting assigns weights based on properties like independence and spread (Brunner et al., 2020). Both  
396 approaches are discussed in this section, complemented by subselection approaches. Model weighting and subselecting to  
397 account for model outliers is discussed specifically in Section 3.3. It is also important to note that some studies may be primarily  
398 interested in assessing the overall performance of the full CMIP ensemble. Applying weighting or model subselection is not  
399 relevant to such analyses.

### 400 **Accounting for Model Performance**

401 Most studies in the literature use simple multi-model means, thus equally weighted MMEs to project future climate change  
402 impacts (Shuaifeng and Xiaodong, 2022). While such approaches capture the overall trends across all models, equal weighting  
403 without any model selection has been criticized for not considering model performance (Shin et al., 2020). Incorporating  
404 information on model skill —by emphasizing better-performing models and down-weighting or excluding models with poor  
405 simulation capabilities— can improve both the accuracy of projections and the assessment of uncertainty in CMIP MMEs  
406 (Merrifield et al., 2020).

407 Weighting, however, is a challenging process as the basis for weights must be determined and other not yet identified but  
408 equally relevant factors may be neglected. Moreover, the relevance of specific model features for a given phenomenon may  
409 change under future climate conditions, making it questionable to assign weights solely based on present-day performance, as  
410 discussed previously in the context of model evaluation. In addition, weighting schemes may inadvertently favor structurally  
411 similar models that produce “mainstream” results, while penalizing outlier models that could provide valuable insights (see  
412 also Subsection 3.3). When bias correction is applied, assessing model performance becomes particularly challenging, as

413 differences between models and observations are largely removed, complicating performance-based weighting. Shin et al.  
414 (2020) addressed this issue by proposing a hybrid weighting approach that preserves performance information while avoiding  
415 unrealistically extreme model weights.

416 Despite these challenges, several studies have demonstrated the potential benefits of performance-based model weighting for  
417 climate projections. Tang et al. (2021) found that weighted MMEs produce more robust projections of extreme precipitation  
418 over the Indo-China Peninsula and southern China than unweighted ensembles. Similarly, Shuaifeng and Xiaodong (2022)  
419 applied a rank-based weighting approach to CMIP6 MMEs for projecting and quantifying uncertainty in cold surges over  
420 northern China. Brunner et al. (2020) discovered a reduction in the projected warming when applying model weighting because  
421 some models showing high future warming have systematically lower performance skills.

422 Another approach to account for model performance is the selection of a subset of models. This can also be considered as a  
423 weighting method, which uses the weight 1 for included models, and the weight 0 for excluded models. MMEs with optimized  
424 sub-selection can reduce the computational load and have been shown to decrease the ensemble-mean RMSE, e.g. by roughly  
425 10–20% for air temperature and approximately 12% for precipitation relative to the full multi-model mean (Hamed et al., 2021;  
426 Herger et al., 2018; Snyder et al., 2024). The central challenges in subselecting are the identification of performance metrics,  
427 as already discussed in Subsection 2.1, as well as the definition of selection criteria that should be made transparent. Herger et  
428 al. (2018) compared different sub-selection approaches, including random ensembles, performance-based ranking, and optimal  
429 ensemble subselection, and found improved performance over the multi-model mean in some cases. In a random ensemble,  
430 multiple models are combined randomly without an explicit optimization strategy. In performance ranking, models are ranked  
431 based on metrics such as accuracy, Q-statistics, mean square error etc. In optimal ensemble sub-selection, a subset of models  
432 is chosen that maximizes performance. Almazroui et al. (2017) similarly found that a subset of the best-performing models  
433 showed better temperature and precipitation projections over the Arabian Peninsula. Numerous further examples exist, e.g.  
434 including the ENSO teleconnection (Roy et al., 2018) and lightning over South/South-east Asia (Chandra et al., 2022).

### 435 **Accounting for Model Dependence**

436 As discussed in Subsection 2.2, climate models are not fully independent. A common approach to address this issue is by  
437 weighting models based on their independence from others. Sanderson et al. (2015) developed a mathematical formulation to  
438 quantify model uniqueness and assign corresponding weights. Boé (2018) argues that model interdependencies are more  
439 effectively assessed through code similarity instead of through result similarity. Although evaluating source code similarity is  
440 indeed challenging (due to issues such as the complexity of model architectures, differing programming languages, licensing  
441 issues and proprietary restrictions), it has the potential to reveal shared model components and algorithms that may not be  
442 evident from model output comparisons alone. Recent model selection approaches also emphasize model independence  
443 (Snyder et al., 2024), with tools such as ClimSIPS explicitly accounting for model dependence (Merrifield et al., 2023).

444 Assessing both source-code similarity and similarity in model results enables the identification of shared methodologies that  
445 can lead to correlated predictions, thereby highlighting potential redundancies within MMEs that may bias ensemble statistics.  
446 Integrating these measures into weighting schemes can therefore improve the robustness of MMEs and contribute to more  
447 reliable and less biased projections.

#### 448 **Combined accounting for Model Performance and Dependence**

449 Knutti et al. (2017) proposed a model weighting method that accounts for model performance as well as model dependency.  
450 This method includes two distance metrics, from models to observations, and among models. Here the “effective repetition of  
451 a model” within an ensemble, outlined by Sanderson et al. (2015), is accounted for, along with the accuracy of a model with  
452 respect to observations.

#### 453 **2.4 Uncertainty Characterization**

454 Model selection and weighting ideally improves uncertainty which remains inevitable when trying to predict climate (Knutti  
455 et al., 2019). Characterizing and understanding it is essential for guiding model evaluation and development, for science and  
456 risk communication, and for assessing climate impacts (Deser et al., 2012a; Deser, 2020; Snyder et al., 2024). When using  
457 future projections from CMIP, three types of uncertainty must be dealt with (Hawkins and Sutton, 2009; Lehner et al., 2020;  
458 Simpson et al., 2021): scenario or forcing uncertainty, natural variability uncertainty, and model uncertainty. The scenario  
459 uncertainty arises because it is not known how human emissions of greenhouse gases and other pollutants from all over the  
460 world will develop in the future, and it is accounted for by modelling different emission scenarios (O’Neill et al., 2014; van  
461 Vuuren et al., 2025). Natural or internal variability uncertainty is due to the chaotic and, thus, unpredictable evolution of the  
462 climate system (Deser et al., 2012b), having a great impact on climate projections (Lehner and Deser, 2023). The unique  
463 realization of our future climate is the response to the combined effect of anthropogenic forcing and internal Earth system  
464 variability. Although internal variability uncertainty cannot be reduced, it is quantifiable (Deser, 2020), and using large  
465 ensembles of a single model is helpful for this purpose (Tebaldi et al., 2021). Finally, the third type –model uncertainty–results  
466 from our imperfect attempts to predict the aforementioned real world realization. It includes differences among models as well  
467 as the varying results that can be obtained within the same model when varying its parameters. While model uncertainty can  
468 be reduced, its interpretation and quantification depend strongly on how the ensemble is constructed (Knutti et al., 2019). An  
469 adequate understanding of uncertainty has the potential to help MMEs users with model selection and thereby reduce  
470 computational burdens (Snyder et al., 2024).

471 Decomposing the total uncertainty of climate estimates into contributions from scenario, internal, and model uncertainty  
472 provides insights into projections’ reliability and potential for reducing uncertainty. This process is called uncertainty  
473 partitioning, and it often involves quantifying the consistency among different members of a MME (Hawkins and Sutton,

474 2009; Lehner et al., 2020; Woldemeskel et al., 2012; Yip et al., 2011). For long-term means of climate data, Hawkins and  
475 Sutton (2009) proposed a widely used method for uncertainty partitioning: they fit a polynomial to each model's output in the  
476 time dimension to separate the forced response from the internal variability. The variance across different model's polynomials  
477 corresponds to the model uncertainty, and the mean of the different residuals across models represents the internal variability.  
478 Finally, the scenario uncertainty is the variance across multi-model means for different forcings. This method assumes (i) that  
479 the forced response can be approximated by the polynomial and (ii) that the arithmetic sum of the different uncertainties  
480 comprises the total uncertainty.

481 To consider the potential non-additive nature of the total uncertainty (ii), Yip et al. (2011) used analysis of variance (ANOVA)–  
482 an approach that partitions the total variance into components due to different sources of variation–to improve the uncertainty  
483 partitioning. Woldemeskel et al. (2012) expanded the uncertainty quantification methodology to include also the spatial  
484 dimension, by introducing the Square Root Error Variance (SREV) method. This method has proven useful for highlighting  
485 regional differences in uncertainty. More recently, exploiting the increasing computational capabilities, Lehner et al. (2020)  
486 overcame the assumption of the polynomial fit (i) from Hawkins and Sutton (2009), which produced significant regional biases,  
487 by using several SMILEs. Instead of calculating the variance of the polynomials as in Hawkins and Sutton (2009), in this  
488 approach, the model uncertainty is calculated as the variance across ensemble means from the available SMILEs. This reduces  
489 methodological assumptions and thereby improves the results, making SMILEs currently a broadly used tool to partition  
490 uncertainty in climate projections. It is important to note that the lack of independence between models (Subsection 2.2), and  
491 the methods to account for it (Subsection 2.3) must also be considered in this context.

492 A question that should be considered, although it can only be partially answered, is whether the MME spread is realistic, too  
493 narrow or too broad. The uncertainty may be too broad if observations are not used correctly to tune models, or if the models  
494 have extensive and diverse structural errors. The ensemble may be too narrow, and thus overly confident if the models are  
495 structurally very similar, if they are overfitted to observations or if uncertain processes are missing. It is also important to  
496 recognize that present-day and future uncertainties arise from different sources: present-day uncertainty mainly reflects the  
497 models' ability to reproduce observations, whereas future uncertainty stems from variations in how models represent physical  
498 processes and feedbacks (Sanderson and Knutti, 2012). Care should be taken when assuming that the spread of present-day or  
499 historical simulations will be the same in the future.

500 As discussed in Subection 2.1, relying solely on how well a model reproduces past climate to assign confidence can be  
501 misleading: models that perform well historically may not accurately project future climate changes, while models that perform  
502 poorly may still provide useful information about future conditions (Hall et al., 2019). An evaluation and uncertainty reduction  
503 technique that avoids this bias is the development of emergent constraints (Hall et al., 2019). An emergent constraint refers to  
504 a statistically robust relationship across a MME between an observable present-day quantity (x) and a projected future change

505 in a quantity ( $\Delta y$ ), typically approximated as linear (Simpson et al., 2021). When this relationship is robust, observations of  $x$   
506 can be used to constrain the plausible range of  $y$ , thereby reducing uncertainty. This is commonly achieved by analyzing the  
507 probability distribution function of  $y$  conditioned on the observed value of  $x$ . This method has been used for assessing the  
508 uncertainty of many processes within different Earth system components (Keenan et al., 2023; Nijssen et al., 2020; Shaw et al.,  
509 2024; Simpson et al., 2021; Smith et al., 2022; Thackeray et al., 2022). ML approaches have also been used to demonstrate a  
510 potential to discover and explore emergent constraints (Nowack et al., 2020). Despite the usefulness of emergent constraints,  
511 care should also be taken when interpreting the results, since the method assumptions may produce overconfident predictions  
512 and may be vulnerable to artifacts within the model (Breul et al., 2023; Sanderson et al., 2021), similar to other uncertainty  
513 reduction methods.

514 While climate models exhibit high confidence in thermodynamic aspects of climate change (e.g. global temperature increase)  
515 due to robust theoretical and observational evidence, dynamic aspects, particularly related to atmospheric circulation, present  
516 significant uncertainties due to their nonlinearity and feedbacks (Shepherd, 2014). Model uncertainties in these two  
517 components are uncorrelated (Zappa and Shepherd, 2017), meaning that errors in one component do not influence or predict  
518 the errors in the other, so separating them allows better understanding of where the biggest uncertainties lie.

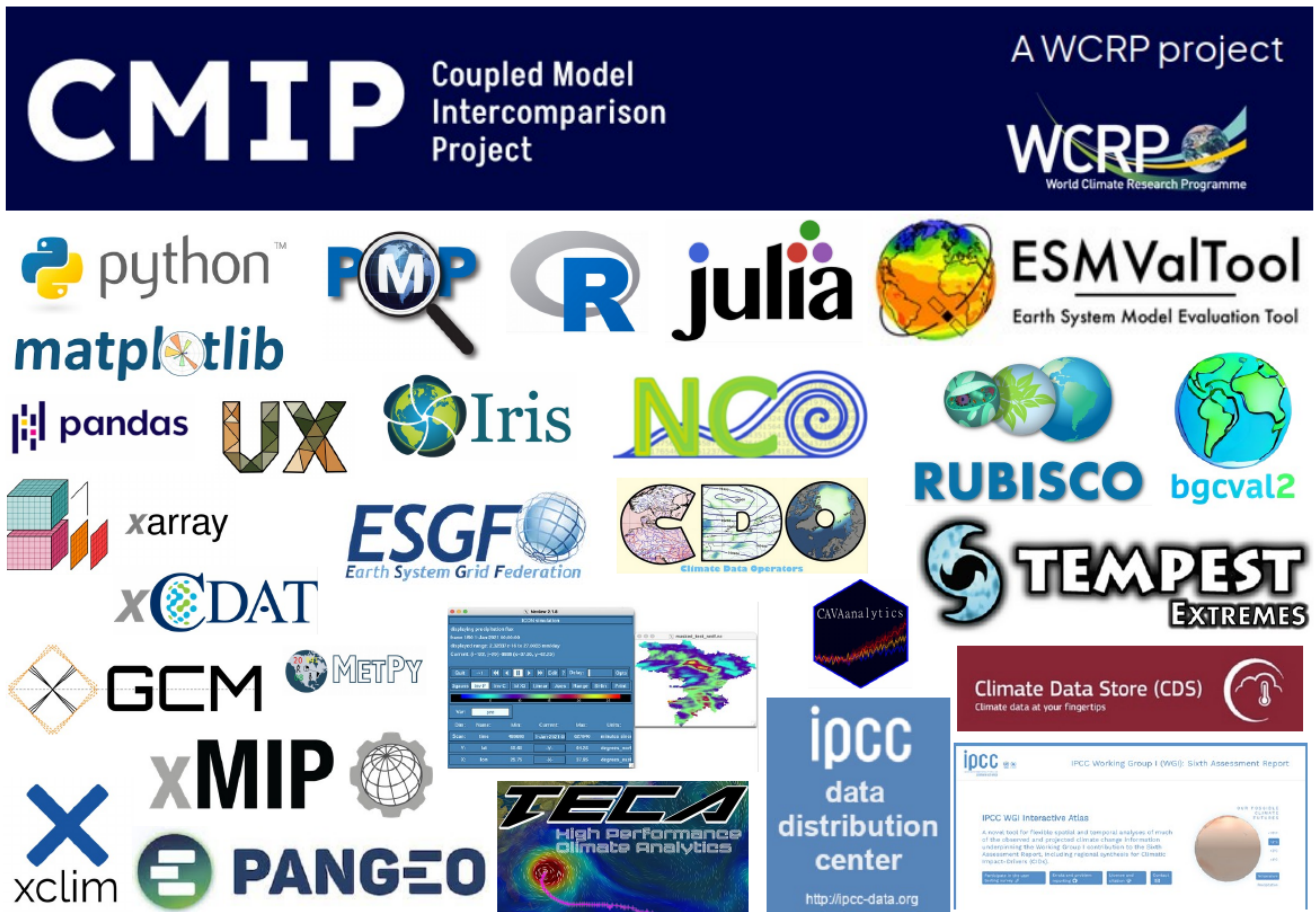
## 519 **2.5 Available Tools for MME Analysis**

520 The analysis of CMIP datasets is greatly facilitated by a variety of tools developed within the global climate science  
521 community. However, the wide range of available tools was not centrally cataloged, making it difficult to obtain a clear  
522 overview of available tools and their capabilities. To address this gap, the WCRP CMIP has undertaken an effort to compile a  
523 central repository (<https://wcrp-cmip.org/tools/>) that encompasses a broad range of resources.

524 The repository includes data access platforms (e.g., Earth System Grid Federation, Climate Data Store, IPCC data distribution  
525 centre, PANGEO, CAVA, Climate Information Portal), which facilitate accessing large and complex data volumes. It also lists  
526 widely used command line operators (e.g., ncview, NCO, CDO) and programming languages suitable for climate data analysis  
527 (such as Python, R, Julia), together with useful packages (e.g., multiple Python packages such as matplotlib, scipy, pandas,  
528 Iris, xarray, xGCM, xMIP, xclim, xCDAT, UXarray, Metpy, aospy). In addition, the repository contains several comprehensive  
529 evaluation and benchmarking tools, such as ESMValTool, bgcval2, RUBISCO, PCMDI Metrics Package, AMBER, and the  
530 MDTF Diagnostic Package. These evaluation tools include diagnostics designed to address specific scientific questions. For  
531 example, ESMValTool incorporates the Climate Variability Diagnostics Package (CVDP, Eyring et al., 2020; Phillips et al.,  
532 2020, 2014; <https://github.com/NCAR/CVDP-ncl>), which facilitates the analysis of modes of climate variability and change in  
533 models and observations (Maher et al., 2024). Another important initiative in process-oriented evaluation is led by the Model  
534 Diagnostics Task Force (MDTF) under NOAA's Climate Program Office (CPO) Modeling, Analysis, Predictions, and  
535 Projections (MAPP) program. It promotes the development and use of process-oriented diagnostics (see Subsection 2.1) in

536 climate and weather prediction models (Maloney et al., 2019; Neelin et al., 2023). Additionally, the WCRP repository includes  
 537 various data analysis and visualization tools, including the IPCC WGI Interactive Atlas, Panoply, TempestExtremes, CAVA,  
 538 TECA, KNMI Climate Explorer, and Google Earth Engine. Figure 3 highlights some of these tools, aiming to promote their  
 539 use across the wider climate community. Basic information about each tool is provided by “Tools description cards” on the  
 540 CMIP website, which include links to tool websites, documentation, tutorials and community support resources. Finally, it  
 541 should be emphasized that the tools repository is actively maintained and continuously updated. To further enhance its utility  
 542 for the broader climate science community, new contributions are highly welcome.

543 While the CMIP tool repository is a key resource for many widely used climate analysis tools, it does not cover all available  
 544 resources. The wider open-source ecosystem - especially within the Python community - offers many additional tools and  
 545 libraries for climate data analysis and is supported by a large and active scientific community on platforms such as GitHub.



<https://wcrp-cmip.org/tools/>

546

547 **Figure 3: Collection of useful tools for using climate data available at <https://wcrp-cmip.org/tools/>.**

### 548 **3. Complementing Topics and Challenges**

549 Building on the general workflow involved in MME studies in Section 2, we draw on the experience within the Fresh Eyes  
550 community to identify common topics and challenges that arise in this context. All of these aspects are also relevant to the  
551 subsections in Section 2; however, our aim here is to provide a dedicated overview of specific topics, allowing researchers to  
552 access the most relevant information in one place.

#### 553 **3.1 Number of Models**

554 Any MME analysis has to face the question of how many models to include, which is not straightforward, as it involves the  
555 trade-offs between model diversity, computational cost, and the accuracy of the results. Increasing the number of ensemble  
556 members has the potential to enhance the robustness of the results by reducing statistical uncertainty. At the same time, state-  
557 of-the-art climate models remain computationally expensive. Downloading and processing these large datasets, particularly in  
558 the context of major intercomparison projects like CMIP, is also a resource-intensive challenge that limits the number of  
559 models included in MME studies. These challenges raise the question how many models are actually required to form a “good”  
560 ensemble size. Here, we focus on the number of models within a MME. A closely related question exists in the context of  
561 large ensembles where the number of perturbed simulations is discussed. Some examples for the number of simulations in  
562 large ensembles are provided in Subsection 4.2. However, many of the arguments and findings in this section apply for both  
563 contexts.

#### 564 **Lower threshold of ensemble size: At least 5 models**

565 If the ensemble size is too small, the inter-model may not be fully captured. This has the potential to lead to an underestimation  
566 of uncertainties and can consequently result in an overestimation of the models’ performance and thus an overconfident  
567 interpretation of the results. It is even possible that a too small ensemble size leads to qualitatively different findings, as shown  
568 by Milinski et al. (2020). In this study, the subsets of two and three models showed a warming after a volcanic eruption, while  
569 the actual known response would be cooling. So, how many models or simulations should be used as a minimum? Several  
570 studies have shown that the error (e.g. root mean squared error when compared to reference data) is reduced substantially up  
571 to about five models in different contexts (Herger et al., 2018; Knutti et al., 2010a; Mendlik and Gobiet, 2016; Milinski et al.,  
572 2020; Steinman et al., 2015). Adding further models is generally beneficial, but the improvement per additional model is much  
573 smaller. Mendlik and Gobiet (2016) find that the subset size can be reduced from 25 to 5 while still being representative for  
574 the entire ensemble. As these studies refer to different quantities and research questions, and were conducted independently,  
575 but still share five as a lower “threshold”, we propose five models/simulations as an initial baseline minimum for MME studies.

576 Depending on the research question however, the minimum number of required models might vary. It can be determined by a  
577 specific method, as explained below.

### 578 **Determining individual minimum ensemble size**

579 If feasible, determining the appropriate minimum ensemble size on a case-by-case basis—depending on the specific research  
580 question and requirements—is preferable to adopting a general minimum. Milinski et al., 2020 proposed a procedure applicable  
581 to diverse research questions. After (1) defining the research question, (2) an error metric (e.g. RMSE) as well as a maximum  
582 acceptable error has to be decided. Then (3), the error for randomly sampled subsets of different sizes has to be quantified.  
583 The number of required models can now be identified as the smallest subset size that has an error below the chosen threshold  
584 (4). If the identified model number is less than half of the initial sample (e.g. the identified subset included 40, thus less than  
585 50 members, when evaluating 100 members) the estimated subset size is robust (5). This requirement is introduced to avoid  
586 resampling bias, as random subsets close to the full ensemble share many members, are no longer independent, and therefore  
587 tend to reproduce the full-ensemble signal by construction rather than providing an unbiased estimate of the required ensemble  
588 size, see Milinski et al., 2020 for details. While this method provides a straight-forward method to identify the ideal minimum  
589 number of models in an ensemble, it requires the availability and analysis of a high number of model simulations.  
590 Consequently, this method might not be feasible for all studies.

### 591 **Remarks for including more models**

592 While the considerations above address the identification of a minimum ensemble size, additional models have the potential  
593 to further improve the model performance. For some applications, larger ensemble sizes are even required, e.g. for the  
594 quantification of internal variability, as estimating higher-order moments of the distribution demands a sufficiently large  
595 ensemble (Milinski et al., 2020). Generally, adding further models improves the statistical robustness of the MME analysis,  
596 but it has to be remembered that the added models should at least partly be independent of the existing models as otherwise  
597 only the weight of single models is increased without any physical reason (Knutti, 2010). See Subsection 2.2 for more details.  
598 A too large ensemble size has also the potential to increase the spread beyond a realistic range as the inclusion of outliers  
599 becomes more probable (Knutti, 2010). Subsection 3.3 therefore discusses strategies for identifying and handling outliers in  
600 more detail. Another consideration becomes relevant when working with different scenarios. As the range of uncertainty  
601 increases with the number of models, using the same number of models across all scenarios is essential to ensure comparability  
602 (Knutti et al., 2010a).

## 603 **3.2 Extremes**

604 Extreme weather and climate events have significant impacts on human society and ecosystems. Understanding the drivers  
605 and producing reliable future projections of these low-frequency high-impact events is therefore essential for effective climate

606 change adaptation planning. When using MMEs to study extreme climate events, ensembles offer both strengths and  
607 challenges.

608 MMEs based on CMIP or CORDEX are widely used both in regional and global studies concerning climate extremes (Kim et  
609 al., 2020; Soares et al., 2023; Vogel et al., 2020; Yang et al., 2012). These studies typically apply statistical approaches, such  
610 as probabilistic modelling, and/or using climate extremes indices defined by the Expert Team on Climate Change Detection  
611 and Indices (ETCCDI). Extreme Value Theory (EVT) provides the theoretical foundation for analyzing extreme events by  
612 offering statistical methods to model the tails of probability distributions (Coles, 2001; DelSole and Tippet, 2022). One widely  
613 used approach within EVT is Generalized Extreme Value (GEV) distribution analysis (Rypkema and Tuljapurkar, 2021), a  
614 statistical framework for modelling the tail of the distribution of rare events. For example, GEV analysis is frequently used to  
615 estimate return periods of extreme rainfall events, allowing assessment of how the frequency and intensity of such events may  
616 change under future climate scenarios (Wehner, 2020). By fitting GEV to observed and modeled data, researchers can evaluate  
617 shifts in extreme event characteristics.

618 A major advantage of using the mean of the MME is that averaging across models reduces noise from internal variability and  
619 thereby amplifies the climate change signal. This can also help identify trends in extreme events (IPCC, 2021). However, the  
620 MME mean might not always be the best choice, particularly when examining the intensity and frequency of extreme events  
621 (Knutti et al., 2010b). Using MME's median or mean can sometimes mask the severity of local extremes, as averaging across  
622 multiple ensemble members can obscure the range of possible outcomes of individual extreme events. This is especially the  
623 case when some models predict significantly different extreme event trends, potentially leading to an underestimation of risks.  
624 Uncertainty remains for both hot and cold extremes, with some models deviating considerably from the multi-model mean.  
625 Uncertainties are particularly large for precipitation extremes, where, despite a general tendency toward heavier precipitation  
626 and longer dry periods, several models project opposing trends in specific regions (Sillmann et al., 2013).

627 In studies on climate extremes, it is therefore important to evaluate how well each model performs in simulating extremes  
628 (Kim et al., 2020; Sillmann et al., 2013) and to correct for biases when appropriate. However, as discussed in Subsection 2.1,  
629 model evaluation is conducted using performance-oriented or process-oriented approaches that generally tend to focus on the  
630 models' ability to capture mean climate states or large-scale circulation patterns rather than models' extreme event  
631 representation. Dedicated evaluations tailored to extremes are therefore required, capturing relevant temporal resolution, region  
632 and variables and comparing to an appropriate reference data set. For example, Kim et al. (2020) evaluated the CMIP6 MME  
633 against ETCCDI climate indices and identified systematic biases, such as a persistent cold bias in cold extremes over high-  
634 latitude regions. When comparing CMIP6 models with CMIP5, they found only limited improvements in simulating  
635 temperature and precipitation extremes, highlighting the need for further advancements in the understanding and representation

636 of extreme climate events in ESMs. As a step following the evaluation, employing model weighting is one possible approach  
637 to address shortcomings, enhancing the accuracy and reliability of extreme event projections (Balhane et al., 2022).

638 When studying extreme climate events, uncertainty is another aspect that is important to account for. As discussed in  
639 Subsection 3.1, ensemble size strongly influences uncertainty estimates, with larger ensembles allowing a more complete  
640 sampling of the range of possible outcomes. In practice, many studies of climate extremes using MMEs rely on a single  
641 ensemble member per model to ensure comparability across models (Kim et al., 2020). However, using only one ensemble  
642 member per model could miss some of the variability in extreme events that larger ensemble runs could capture, particularly  
643 as often not too extreme members are submitted for intercomparison projects as CMIP. Nevertheless, given the constraints on  
644 computational resources and the availability of large ensembles, this method remains a common compromise.

645 While increasing ensemble size can help reduce uncertainties, it does not eliminate the limitations inherent to individual  
646 models. Downscaling techniques, either statistical or dynamical using RCMs, can provide higher-resolution data to improve  
647 the representation of extremes in specific regions. For example, the bias-adjusted high-resolution RCM outputs in the EURO-  
648 CORDEX project showed an improvement in the simulation of extreme temperature and precipitation indices across Europe,  
649 underscoring the value of RCMs for more reliable and region-specific climate projections (Coppola et al., 2021; Dosio, 2016).  
650 MMEs based on RCM projections are particularly valuable for highly vulnerable regions, offering insights into potential  
651 changes in local extreme events (Dosio, 2017; Tegegne et al., 2021) and supporting planning for challenges such as water  
652 scarcity, food security, and disaster preparedness. For more details regarding downscaling, see also Subsection 3.4.

### 653 **3.3 Outliers**

654 Outlier models have at times been disregarded in MME analyses because convergence toward the ensemble mean has been  
655 interpreted as a measure of model reliability. However, the use of convergence as a measure of model reliability has been  
656 criticized because it favors simulations that are closer to the multi-model mean, while underrepresenting uncertainty across a  
657 wider range of plausible outcomes (Tebaldi and Knutti, 2007). For example, the original version of the reliability ensemble  
658 average (REA) weighting method assigns higher weights to models that better reproduce the current climate, but also penalizes  
659 models that diverge from the ensemble mean (Giorgi and Mearns, 2002). As a result, outliers receive lower weights even when  
660 their differences may reflect physically plausible behavior rather than poor model performance. This is especially problematic  
661 because convergence toward the ensemble mean can partly arise from genealogical similarities among models, rather than  
662 independent confirmation of a result (Tebaldi and Knutti, 2007). Despite these concerns, there is a history of privileging  
663 convergence towards the MME mean within the climate science community. For example, in the third IPCC assessment report  
664 two models were discarded because of extreme warming projections associated with very high climate sensitivity (Tebaldi and  
665 Knutti, 2007). More recently, convergence-based ideas continue to influence MME subsetting approaches (Palmer et al., 2023)  
666 and are still used, at least in part, in MME evaluation frameworks (Amali et al., 2024). Whether emphasizing convergence is

667 appropriate depends on the purpose of a given study. For applications such as the analysis of climate extremes, averaging  
668 across models can mask the full range of possible outcomes. In these cases, outlier models may provide valuable information  
669 about plausible high-impact scenarios rather than representing spurious deviations from the ensemble mean. See also  
670 Subsection 3.2 for more details on extremes in MMEs. Building on the insights to weighting and building subsets of models  
671 in Subsection 2.3, we discuss here in more detail how and when to account for outliers.

## 672 **Exclusion of Outliers**

673 One approach to account for outliers, is exclusion, meaning the removal of models with outlier status from an ensemble. While  
674 this can help reduce unrealistic spread and improve agreement with observations, it carries the risk of omitting simulations  
675 that represent rare but physically plausible events. In some cases, however, the benefits of exclusion outweigh the drawbacks.  
676 For example, Mudryk et al. (2020) identified outlier models for some seasons and regions in their study of snow cover change  
677 in the Northern hemisphere. These models, which overestimated snow cover in areas of low snow mass, were excluded to  
678 improve the alignment between observational data and CMIP6 MME projections. Similarly, the Swiss Climate Scenarios  
679 CH2018, based on EURO-CORDEX, excluded some outlier GCMs to narrow uncertainty ranges for temperature and  
680 precipitation (Sørland et al., 2020). The consequences of outlier inclusion or exclusion have been explored in the literature.  
681 Sun and Archibald (2021) compared “aggressive” and “conservative” approaches—respectively including and excluding  
682 outliers—and found that, for their study, the differences in results were relatively minor. Similarly, Bracegirdle and Stephenson  
683 (2012) presented analyses both with and without outliers to illustrate the sensitivity of polar warming estimates to outlier  
684 inclusion and different forms of regression. Overall, while models with outlier projections may be excluded to improve MME  
685 alignment with observations or to reduce uncertainty, this should be done with caution. Exclusion is most justified when a  
686 model is known to be deeply flawed. Otherwise, removing projections of rare but plausible events may limit the assessment  
687 of adaptation strategies and risk management options (Knutti et al. 2010).

## 688 **Weighting or Penalization**

689 MME inclusion of outlier models is often accomplished through weighting. One commonly used approach is weighting based  
690 on root-mean-square error (RMSE) skill scores. For example, Tegegne et al. (2020) preserve MME spread by applying the  
691 Katsavounidis–Kuo–Zhang algorithm to select ensemble members based on their contribution to representing the full range of  
692 variability within the sampling space for extreme indices of interest, as recommended by World Meteorological Organization’s  
693 ETCCDI. The IPCC characterizes this approach as suitable for detecting “moderate extremes”—events expected to occur up to  
694 10% of the time (Seneviratne et al., 2012). In this approach, detecting extremes prior to taking the MME mean is useful for  
695 weighting members such that the full range of variability within the MME is largely preserved. To identify and characterize  
696 more extreme events, methods based on extreme value theory (EVT) are required. EVT focuses on values located in the very

697 ends of tails of probability distribution functions (PDFs) and is therefore better suited to representing rare, high-impact  
698 extremes (DelSole and Tippett, 2022).

699 Among weighting methods rooted in classical statistics is the use of outlier insensitive methods. These are methods that retain  
700 outlier models while limiting their influence on the ensemble result. Such methods include using the ensemble median instead  
701 of its mean as a measure of the MME's center which reduces sensitivity to outliers (Ge et al., 2021). Rank based tests of  
702 statistical significance provide another option. Because they rely on the rank, or position of a value within a PDF, rather than  
703 the value of a particular data point within a sample, these are largely insensitive to outliers (DelSole and Tippett, 2022). This  
704 test is recommended by the World Meteorological Organization for hydrological data analysis and is robust to outliers and  
705 non-normal data distributions (Rojpratak and Supharatid, 2022).

706 Weighting methods based on Bayesian statistics have been developed to sample uncertainty across a broad statistical space.  
707 Compared to the frequentist statistics which uses a fixed population parameter to describe probability distributions, Bayesian  
708 statistics uses a conditional parameter that depends on the shape of the PDF for a given dataset (Clyde et al., 2022). Xu et al.  
709 (2019) apply Bayesian model weighting to statistically downscale precipitation data for site-specific analyses. The authors  
710 argue for the use of statistical downscaling due to its relatively low computational expense with finer spatial and temporal  
711 resolution data. At the same time, they note that dynamic downscaling can underestimate extremes and be overly sensitive to  
712 outliers. These limitations are addressed by applying a Bayesian weighted average, which reduces outlier influence while  
713 retaining information about uncertainty.

714 Penalization refers to methods that are explicitly designed to reduce the weight of outlier models within an MME. This can be  
715 achieved through bias correction and ridge regularization. In Shin et al. (2020) outlier models are defined as those that generate  
716 projections that are unusually close to the hydrological variable in observation data, which the authors attribute to excessive  
717 regional calibration to observations. To limit the influence of such models, they propose a hybrid method combining Bayesian  
718 weighting with bias correction. Ridge regularization, frequently used in ML context, is a form of linear regression that  
719 incorporates a penalty term to constrain variables with unusually strong linear correlations, thereby reducing the risk of  
720 overfitting. Labe and Barnes (2022) apply ridge regularization to limit the sensitivity of an artificial neural network to outlier  
721 influence.

### 722 **3.4 Downscaling Techniques**

723 Acquiring regional information on climate change is essential for impact, vulnerability, and adaptation studies. While CMIP  
724 GCMs are internationally established sources for climate projection data, their typical grid resolution of 100–250 km (Liang-  
725 Liang et al., 2022; Weigel et al., 2010) limits the ability to provide locally relevant information (Grose et al., 2023).  
726 Downscaling is therefore necessary, especially for regions with complex topography or localized climate phenomena (Wilby

727 and Fowler, 2010). Various downscaling techniques exist, including statistical downscaling (Gebrechorkos et al., 2023;  
728 Wootten et al., 2024), dynamical downscaling (Knutson et al., 2013; Tapiador et al., 2020), and novel machine-learning based  
729 approaches (Sachindra et al., 2018; Soares et al., 2024), each with its own strengths and limitations (Hall, 2014). These  
730 downscaling approaches differ from the regriding methods discussed in Subsection 3.5, which are purely mathematical  
731 interpolation techniques and do not introduce additional physical or statistical information.

732 The statistical downscaling technique uses statistical relations between coarse-resolution GCM climate data and observed local  
733 climate data to generate fine-scale downscaled projections for a specific region (Oxarart and Parker, 2024). Its reliability  
734 depends on the quality of observational data and on the assumption that calibrated relationships remain valid in a changing  
735 climate. Statistical downscaling is computationally efficient and can reduce the cool bias compared to the original CMIP  
736 simulations, as shown e.g. by Xu and Wang (2019) applying the Bias Correction and Spatial Downscaling (BCSD) technique  
737 for daily maximum temperature over China. However, it may struggle to represent non-linear processes or unprecedented  
738 extremes if these are not well captured in the historical record.

739 In dynamical downscaling, output from a GCM is used as boundary conditions for a RCM, which simulates climate on a  
740 limited-area domain and hence employs a finer resolution (Di Luca et al., 2015). Specifically, the CORDEX (Giorgi, 2019;  
741 Gutowski Jr. et al., 2016) initiative provides an international framework in which multiple institutions generate and evaluate  
742 such regional climate projections driven by GCM simulations. Dynamical downscaling can capture regional physical processes  
743 that GCMs cannot resolve (Giorgi and Gutowski, 2015). However, it is computationally demanding and depends on the  
744 availability of suitable RCMs. In addition, systematic biases in the GCM can propagate into, and potentially degrade, the  
745 regional simulations (Di Virgilio et al. 2022). One prominent example of the application of dynamic downscaling is the  
746 derivation of European Centre for Medium-Range Weather Forecasts Reanalysis 5 (ERA5) dataset (Xu et al., 2021). Another  
747 example for the Australian region is Grose et al. (2023) who used CMIP6 multimodel ensemble downscaling to provide  
748 accurate, scenario-based climate change projections. Liu et al. (2021) found that dynamical downscaling does not necessarily  
749 perform better compared to dynamic downscaling approaches.

750 ML-based downscaling methods have recently been applied to generate high-resolution projections from GCM simulations,  
751 leveraging their ability to capture complex and non-linear relationships. They can also efficiently process large datasets and  
752 integrate multiple variables, which has potential to improve downscaling results (Rampal et al., 2024). See Subsection 4.1 for  
753 details and examples.

### 754 **3.5 Regriding Techniques**

755 Each model produces output on its own underlying grid, often referred to as the ‘native’ grid. When combining models with  
756 different native grids into a MME, researchers must decide on whether to keep the native grids or to regrid their data to a

757 uniform grid. One option is to retain native grids and avoid regridding altogether. Analysing and visualization individual MME  
 758 model results in the models’s native grid is one way to accomplish this (Quesada et al., 2017). Another approach is to compute  
 759 zonal means for each model and then average these, which allows results to be combined without regridding (Boysen, 2020).  
 760 In practice, however, most MME studies regrid model output to a common grid to ensure spatial consistency prior to analysis.  
 761 Regridding introduces several methodological choices related to both spatial and temporal dimensions, including the selection  
 762 of a target grid resolution (coarser, intermediate, or finer) and type, the interpolation method, and the treatment of differing  
 763 model calendars.

764 Let us consider the question of which grid resolution to choose. A range of grid resolutions are likely to exist within an MME,  
 765 with one or more of those grids being at the coarse end of the range. Regridding to a coarser grid can improve computational  
 766 efficiency, and facilitate comparison across models, but may smooth spatial gradients and dampen localized extremes.  
 767 Conversely, regridding to a finer grid can better preserve small-scale features and extremes, but does not add new physical  
 768 information and may give a false impression of increased spatial detail. Many studies that mention regridding do not explain  
 769 the direction of regridding or the rationale behind it (Achugbu et al., 2022; Cook et al., 2020; Gergel et al., 2024; Hong et al.,  
 770 2022; Song et al., 2021; Zhao and Dai, 2021), showing that it is common in literature to not disclose the details of regridding.  
 771 Where documentation on details is provided, different motivations are evident. Teuling et al. (2019) regrid to a coarser grid  
 772 only for data visualization purposes. Iles et al. (2020) state that regridding to a finer grid has the ability to preserve localized  
 773 extremes to a greater degree than lower resolution data.

774 Next, one must consider how to interpolate the data that is being regridded. In many Python-based workflows, the default  
 775 interpolation method is bilinear interpolation. This is suitable for many, but not all, variables depending on the type of analysis  
 776 that is being carried out. Table 1 provides an overview of the available interpolation methods, which data types they should be  
 777 applied to, and some examples of CMIP variables for each data type. Beyond the presented methods, additional interpolation  
 778 methods exist (National Center for Atmospheric Research Staff 2014).

779 **Table 1.** Interpolation methods commonly used in climate data analysis

Interpolation method	When to use	Data type	Example variables
None	When no filling or averaging of the original data is desired	Categorical	treeFrac, cropFrac
Bilinear	When data point values vary smoothly across a surface	Continuous	tas, sst

First-order conservative	When fluxes must be conserved over a given area	Conservative	pr, evspsbl
Second-order conservative	When fluxes must be conserved over a given area (smoother than first-order conservative when going from coarser to finer grid)	Conservative	mrro, mrso
Nearest neighbor	When strong contrast between areas with discrete or categorical values must be maintained	Categorical	treeFrac, cropFrac
Patch	When the computation of accurate derivatives is needed	Conservative	tauu, tauv

780

781 In addition to the variety of spatial resolutions present within an MME, temporal inconsistencies may also exist among  
782 members. These arise because models can use different calendar conventions when storing output in the commonly used  
783 netCDF format, which supports nearly ten calendar types (NetCDF Users Guide: NetCDF Utilities, 2025) and subsequent  
784 different encoded calendars in the modelling centers. The best choice of calendar for a given study will depend on the study  
785 particulars and researcher preference. Regardless of this choice, calendars should be brought into alignment during the  
786 regridding process to avoid inconsistencies in subsequent MME analyses.

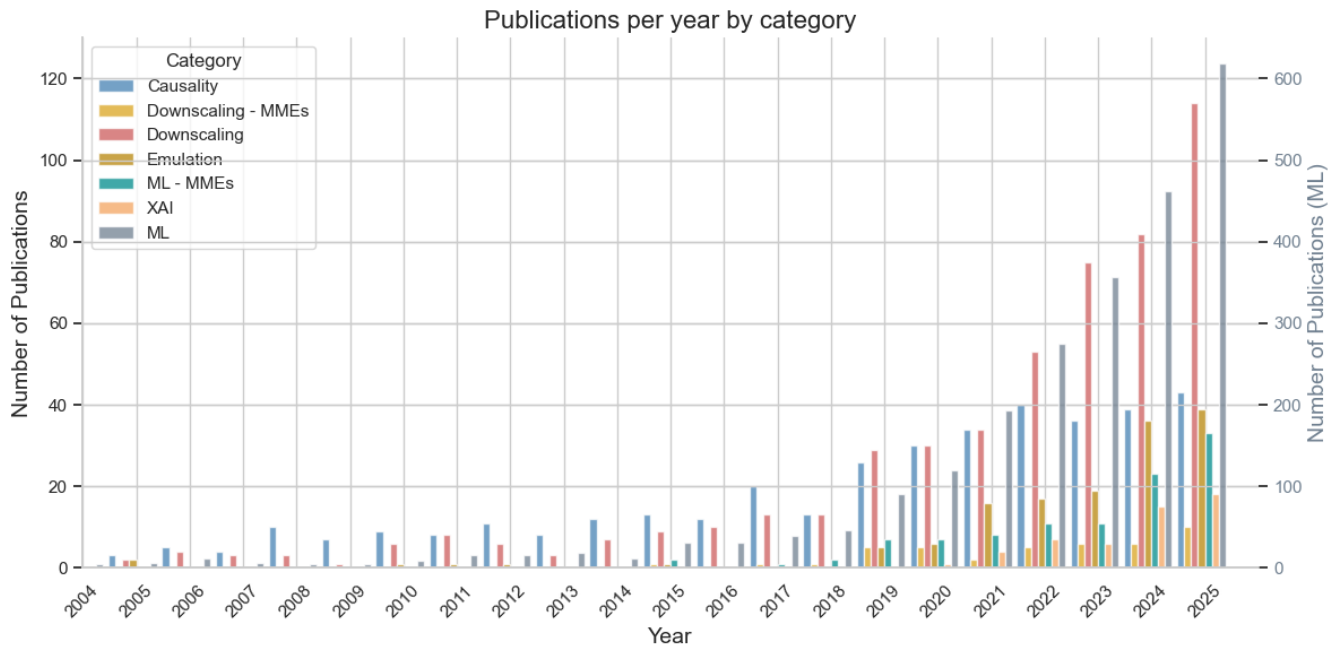
## 787 4. Outlook

### 788 4.1 Machine Learning

789 With the rapid production and accumulation of climate data, automated and increasingly sophisticated analysis techniques  
790 have become essential (Rupe et al., 2017; Glymour et al., 2019). Recent advances in ML have been most pronounced at weather  
791 timescales, where large training datasets enable robust data-driven approaches. In contrast, applications to climate timescales  
792 are more challenging due to limited sample sizes, stronger non-stationarity, and a frequent occurrence of conditions outside  
793 the range of the training data, especially for extreme events. Despite these challenges, ML has emerged as a valuable tool for  
794 enhancing ensemble approaches in climate science (see Fig. 4).

795 Over the past 5-10 years, ML applications have demonstrated significant advantages in addressing non-linear, high-  
 796 dimensional, and hierarchical problems (Li et al., 2021 and references therein). By utilizing observational data as either a  
 797 reference, benchmark, or a constraint, ML offers significant potential to extract additional insights from MMEs. Approaches  
 798 based on neural networks, causal inference, explainable artificial intelligence (XAI), and nonlinear multivariate emergent  
 799 constraints have become increasingly competitive with traditional numerical, knowledge-based methods (see Fig. 5 and de  
 800 Burgh-Day and Leeuwenburg, 2023; Eyring et al., 2024). These capabilities make ML particularly well-suited for identifying  
 801 patterns and complex physical processes within climate model data, enabling a more comprehensive exploration of the valuable  
 802 information embedded within the data (Reichstein et al., 2019; Wang et al., 2018). In short, ML has the potential to make  
 803 climate models better and faster, while reducing their high energy consumption. Nevertheless, the application of ML  
 804 algorithms in constructing MMEs for climate impact assessments remains at an early stage. Below, we provide an overview  
 805 of emerging ML approaches for analyzing MMEs.

806



807

808 Figure 4. Number of publications per year involving different ML-related techniques in the context of climate modelling: ML  
 809 (y-axis on the right), ML and MMEs, Downscaling, Downscaling and MMEs, Causality, Emulators, and Explainable AI (XAI).  
 810 The data was extracted from the citation reports available at Web of Science  
 811 (<https://www.webofscience.com/wos/woscc/basic-search>) using the queries provided in Appendix A.

## 812 **Downscaling and Bias Correction**

813 ESMs have horizontal resolutions often far coarser than those needed by decision makers and also exhibit substantial biases  
814 (Maraun et al., 2017). Recently, ML has been exploited to bias-correct and downscale MME’s outputs—with both processes  
815 often done simultaneously. The source data are typically coarse-resolution outputs from climate models for historical periods,  
816 while the targets are high-resolution observational datasets, such as gridded products interpolated from gauge stations. It is a  
817 common practice to perform dimensionality reduction (e.g. performing principal component analysis on the raw data) before  
818 the training process. Multiple ML methods have been tested and compared, including random forests, support vector machines,  
819 relevance vector machines, and artificial neural networks (Crawford et al., 2019; Sachindra et al., 2018; Wang et al., 2018; Xu  
820 et al., 2020; Dey et al., 2022; Jose et al., 2022; Shetty et al., 2023; Zebarjadian et al., 2024; Li et al., 2021). The domain of  
821 these studies is generally limited to river basin scales (Crawford et al., 2019; Dey et al., 2022; Jose et al., 2022; Sachindra et  
822 al., 2018; Shetty et al., 2023; Xu et al., 2020; Zebarjadian et al., 2024), although Wang et al. (2018) and Li et al. (2021) obtained  
823 good results at a country level for Australia and China. In many of these studies, it has been found that tree-based approaches  
824 (like random forests) commonly perform better than other algorithms, making them a strong baseline for future research that  
825 aims to improve bias correction or downscaling algorithms.

826 Although these approaches provide a practical way to leverage MME future projections and observations to obtain a “best  
827 estimate” of future quantities, there are several critical limitations to consider. First, within these methods, it is assumed that  
828 the relationships between model outputs and observations remain stationary, including model biases and errors (Maraun, 2016).  
829 However, skillful or poor model performance during the historical period does not necessarily translate into the same for the  
830 future. Training the algorithms with historical data can thus lead to projections becoming overly constrained. Potential  
831 solutions for this aspect are to use trend-preserving learning (Wang and Tian, 2024) or climate-invariant ML methods (Beucler  
832 et al., 2024).

833 Another critical aspect that requires further attention in future ML-based bias correction and downscaling efforts is the potential  
834 degradation of the representation of temporal variability in final estimates (Shetty et al., 2023). Among the studies mentioned  
835 above, only Li et al. (2021) acknowledged that their model outputs showed a significant reduction in the interannual variability  
836 relative to the original CMIP models. Thus, it is necessary to implement evaluation metrics for the algorithms that consider  
837 aspects such as the standard deviation of the generated time series, the frequency and persistence of extreme events, and the  
838 amplitude of different modes of variability. Most approaches aim to minimise only one error metric ignoring the skill regarding  
839 other aspects and the physics behind them. For example, the mean precipitation could be improved but the representation of  
840 the extreme events or the number of wet days may not be addressed. Algorithms that can minimize multiple loss functions  
841 simultaneously could be advantageous to preserve multiple statistical features of the fields of interest (Lin et al., 2019; Sener  
842 and Koltun, 2018; Zuluaga et al., 2013). Furthermore, ML-based approaches normally focus on predicting just one variable.

843 Using methods that aim to predict multiple variables could help preserve inter-variable relationships (while also helping  
844 preserve different modes of variability).

845 Finally, most bias correction or downscaling algorithms are trained to predict the outputs in one grid cell based on the nearest  
846 CMIP grid cell. This approach neglects spatial relationships contained either within the inputs or the desired outputs. ML  
847 methods that account for spatial relationships could be of use, including convolutional neural networks (Gu et al., 2018; LeCun  
848 et al., 2015; Wang and Tian, 2022, 2024). Incorporating spatial relationships, multiple variables, and multiple error metrics,  
849 also diminishes the impact of observational uncertainty, since physical relationships are more easily preserved, and it also  
850 reduces the risk of producing overly constrained projections. Considering the limitations of the approaches mentioned for  
851 detecting physically plausible connections, it is essential to explore additional methodologies, with causal inference being one  
852 promising option.

### 853 **Causal inference for climate models**

854 Causal inference aims to identify the causal structure of complex systems such as the Earth and to quantify causal effects by  
855 combining domain knowledge, ML models, and data from observations and climate model simulations (Runge et al., 2023,  
856 and references therein). Because widely adopted methods based on simple descriptive statistics often fail to accurately capture  
857 the underlying physical mechanisms (Beven and Freer, 2001), structural causal models (SCMs) have become a well-  
858 established approach in statistics and ML for causal inference (Runge et al., 2019). An important resource in this area is the  
859 causality benchmark platform **causeme.net**, which aims to provide benchmarks with well-established causal structures (Runge  
860 et al., 2020). Process-oriented causal analysis has been applied across a broad range of topics, including Arctic processes and  
861 their connections to the mid-latitudes (Docquier et al., 2022, 2024; Galytska et al., 2023; Kaufman et al., 2024; Kretschmer et  
862 al., 2020; Polkova\* et al., 2021), Atlantic–Pacific interactions (Karmouche et al., 2023) and subpolar gyre variability (Falkena  
863 and von der Heydt, 2024).

864 Building on this foundation, recent research has increasingly explored the integration of causal discovery with deep learning  
865 (DL), presenting a promising avenue for improving climate simulations (Iglesias-Suarez et al., 2024; Kyono et al., 2020; Luo  
866 et al., 2020; Russo and Toni, 2022; Wang et al., 2024; Yoon and Schaar, 2017; Zhang et al., 2023). This combination aims to  
867 address biases and uncertainties associated with subgrid-scale processes, such as clouds and convection. Previous research has  
868 demonstrated DL's capability to effectively represent small-scale processes, such as deep convection, using storm-resolving  
869 model simulations (Eyring et al., 2021; Gentine et al., 2018; Grundner et al., 2022). Despite this potential, DL algorithms have  
870 been criticised for robustness issues, poor generalization, and the reliance on spurious, non-physical relationships, particularly  
871 when conditions diverge from the training data (Brenowitz et al., 2020; Scholkopf et al., 2021; Thuy and Benoit, 2024).  
872 However, Iglesias-Suarez et al. (2024) demonstrated that causal discovery can effectively identify the physical drivers of  
873 subgrid-scale processes, thereby enhancing the reliability of DL algorithms. Their causally-informed approach generates

874 climate means and variability that closely match original simulations, while preventing spurious links typically seen in  
875 traditional DL-based parameterizations. This aligns with previous work by Zhang et al. (2023) emphasizing the value of  
876 integrating domain knowledge to address the limitations of purely data-driven models. While these studies currently do not  
877 pertain directly to multi-model analysis, their methodologies hold significant potential for future applications in this area.

### 878 **Process-oriented causal model evaluation**

879 Building on the identification of causal relationships, these frameworks also offer opportunities for systematically evaluating  
880 climate models. Detecting similar causal connections in observations and model simulations provides an opportunity to assess  
881 model performance that indicates whether models can correctly reproduce processes in the climate system. Such an evaluation  
882 framework was first introduced by Nowack et al. (2020) and was termed causal model evaluation (CME). To facilitate the  
883 comparison of causal relationships, the authors introduced a modified asymmetric  $F_1$  score metric to classify the agreement  
884 between compared causal graphs. A similar approach was proposed by Vázquez-Patiño et al. (2020). Debeire et al. (2025)  
885 built their study upon the findings of Nowack et al. (2020) to address the practical challenges of integrating CME for CMIP6  
886 MMEs projections. The authors adopted and adjusted the  $F_1$  **score definition and complemented it with a**  
887 **distance metric  $1 - F_1$**  with smaller distance values indicating greater similarity, both in terms of performance relative  
888 to the reference graph and in terms of dependence among the models. Based on this metrics, Debeire et al. (2025) developed  
889 a new weighting scheme, termed causal weighting, inspired by the earlier works of Knutti et al. (2017) and Brunner et al.  
890 (2020), which accounts for both model performance and interdependence of causal networks.

891 Ricard et al. (2024) employed a network-based approach, termed netCS, which leverages sea surface temperature (SST)  
892 variability and teleconnections to constrain Equilibrium Climate Sensitivity (ECS) and Transient Climate Response (TCR).  
893 The authors argue that the behavior of SST networks serves as a reliable proxy for how models respond to increased  $CO_2$   
894 concentrations. Their results show that some models capture regional SST distributions well but fail to replicate connectivity  
895 patterns, and vice versa. This distinction is crucial for evaluating model performance over historical periods, as models that  
896 realistically reproduce past SST patterns may exhibit more physically consistent behavior, even if they are not necessarily  
897 better tuned. While this does not guarantee that those models are superior for future projections (Rasp et al., 2018; Zhu and  
898 Poulsen, 2021), it provides valuable evidence. The authors further propose that causal networks, when used alongside  
899 traditional emergent constraints, offer a more reliable framework for ranking climate models in future climate projections.

### 900 **Machine Learning for Climate System Emulation**

901 Climate model emulators, including surrogate models, are simplified representations of the complex processes embedded in  
902 climate models, enabling faster computations and predictions. They mimic the behaviour of a climate model without explicitly  
903 solving the underlying equations. ML presents a unique opportunity to emulate components of the climate system through

904 novel and computationally efficient parameterizations. Such approaches have the potential to increase the efficiency of climate  
905 simulations while enabling higher resolution simulations (Eyring et al., 2021; Gentine et al., 2018). However, the success of  
906 ML emulation of the climate system varies depending on the choice of algorithm, temporal resolution, type of training data,  
907 and model complexity (Dueben and Bauer, 2018; Scher, 2018).

908 One notable initiative in this context is ClimSim, a hybrid physics-ML dataset designed to provide high-quality data for training  
909 ML emulators of climate processes (Yu et al., 2023). These datasets have been tested for deterministic and stochastic  
910 parameters, and show promise for future climate simulations if applied properly. Future studies could explore the use of MMEs  
911 as training data to develop novel ML-based emulators. Complementing the available data to train emulators, Lu and Ricciuto  
912 (2019) demonstrate an innovative approach integrating SVD, Bayesian optimization, and neural networks to create a  
913 computationally efficient surrogate model. Weber et al. (2020) provide technical notes of ML, using the example of forecasting  
914 precipitation under CO<sub>2</sub> forcing. The remarkable computational efficiency and ability of ML emulators to replicate complex  
915 climate processes with high precision demonstrates their potential. Nevertheless, several challenges remain, including the high  
916 computational cost of running ML models, limited diversity in training data, and the need for more robust methods to evaluate  
917 simulations.

### 918 **Further promising Future ML Avenues**

919 There are numerous promising avenues involving ML for the analysis and processing of CMIP outputs. Explainable AI (XAI),  
920 which aims to extract physical insight from otherwise black-box ML models, offers substantial potential for identifying  
921 physically meaningful changes in the Earth system simulated by CMIP models (Rader et al., 2022). For example, layer-wise  
922 relevance propagation (LRP), has been used to identify the spatial regions and input features that neural networks rely on when  
923 generating predictions (Toms et al., 2020) and has proven especially valuable for improving interpretability by visualizing  
924 relevance heatmaps (Hilburn et al., 2020; Labe et al., 2024; Labe and Barnes, 2022; Sonnewald and Lguensat, 2021). This  
925 interpretability adds value to ensemble evaluation, providing critical information that can inform model weighting schemes.  
926 ML also offers effective tools for evaluating both the performance and independence of climate models within MMEs, further  
927 supporting the development of ensemble weighting metrics (Brunner and Sippel, 2023). These types of methods also support  
928 the development of process-oriented bias correction and downscaling methods for MMEs (Maraun et al., 2017). Furthermore,  
929 efforts to predict end-user-relevant variables that are not directly simulated by GCMs, including crop yield (Crane-Droesch,  
930 2018; Sidhu et al., 2023; Veenadhari et al., 2014) and power generation potential (Jung et al., 2021; Nwokolo et al., 2023;  
931 Yeganeh-Bakhtiary et al., 2022), demonstrate the potential of ML to enhance the applicability of MMEs for stakeholders and  
932 decision-makers. Given ML's growing role in improving climate projections, interpretability, and practical usability, AI-ready  
933 databases such as ClimateSet (Kaltenborn et al., 2023) represent valuable resources for the research community.

#### 934 **4.2 Single Model Initial Condition Large Ensembles (SMILEs)**

935 For some activities and experiments, many models in CMIP5 and CMIP6 provide only one ensemble member (Milinski et al.,  
936 2020; Olonscheck and Notz, 2017). Consequently, modelling groups strive to provide their best performing members, carefully  
937 calibrated to the same internationally available observational datasets. This implicit incentive for modelling groups to add  
938 simulations to the CMIP MME that are less extreme can lead to a MME that underestimates the uncertainties. The outcome  
939 of this incentive is shown by findings from Sanderson et al. (2008) who found that the standard model performed comparatively  
940 to the best-performing model. To address this and other limitations of MMEs based on single-members from different  
941 modelling centers, an emerging approach (see Figure 5) is to include multiple simulations from individual models that differ  
942 only in their initial conditions (Maher et al., 2021). When there are at least 10 ensemble members, we refer to the ensembles  
943 as SMILE (Deser et al., 2020).

944 Although MMEs are useful for examining the combined influence of three types of uncertainties in climate projections (model  
945 uncertainty, internal variability uncertainty, and scenario uncertainty), MMEs do not allow us to distinguish internal variability  
946 from the forced response. On the other hand, SMILEs allow us to quantify internal variability for any given future scenario,  
947 independent from model uncertainty (Deser et al., 2012b; Lehner et al., 2020). This capability is particularly powerful for  
948 regional detection and attribution studies and for the analysis of extreme climate events (Lehner et al., 2017; McKenna and  
949 Maycock, 2021; von Trentini et al., 2020; van der Wiel et al., 2021; Pérez-Carrasquilla et al., 2025).

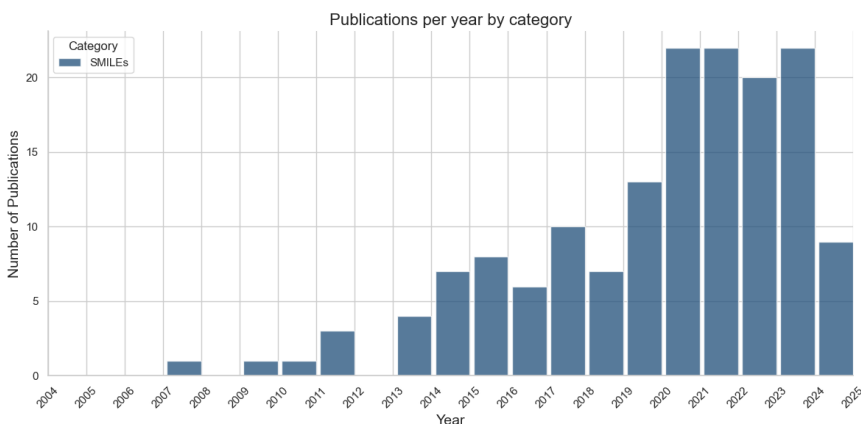
950 While large ensemble simulations are well established as essential for studying extreme events (see also Subsection 3.2), they  
951 are equally relevant for the analysis of compound events that result from combinations of multiple weather and climate drivers  
952 (e.g. simultaneous drought and heatwave; Bevacqua et al., 2022, 2023; Wu et al., 2023). Such events are characterized by  
953 complex interactions between extreme conditions across variables, space, or time. Consequently, single variable based  
954 approaches may underestimate risks, and neglect the impact of the interplay between different variables. In this context,  
955 Bevacqua et al. (2023) showed that attributing compound events requires larger sample sizes than events based on extremes  
956 in a single variable, especially when the drivers are weakly correlated and have similar trends. Sampling a wide range of  
957 possible atmospheric conditions using SMILEs helps avoid underestimating the frequency and severity of compound events  
958 and provides deeper insights into their physical drivers and potential future changes.

959 When multiple realizations are available, it is considered good practice to average all realizations from each model and  
960 incorporate these means into the MME. This approach is appropriate for applications focusing on long-term mean changes or  
961 the forced response, but is not appropriate for analyses of extremes or internal variability, where individual realizations carry  
962 relevant information. When more ensemble members are used, it is important to remember that the ensemble size available for  
963 the individual models should not influence the weight given to this model in the MME (Knutti et al., 2010a). Future studies

964 should provide a methodological framework on how to combine SMILEs and MMEs in the most productive and meaningful  
965 way.

966 Some modelling centers have also applied SMILEs framework to partition the forced response into contributions from  
967 individual forcings (e.g. greenhouse gases, aerosols, biomass burning) with single-forcing large ensembles (SFLE; e.g. derived  
968 from CESM2 framework; Simpson et al., 2023). These ensembles change only one forcing at a time while holding all others  
969 fixed, thereby enabling attribution of the drivers underlying responses identified in all-forcing SMILEs.

970 One challenge to employing SMILEs is data accessibility. To address this, the Multi-Model Large Ensemble Archive  
971 (MMLEA) was developed (Deser et al., 2020). The newly published MMLEAv2 expands upon the original archive by  
972 including a larger number of models (18 compared to 7 previously) and more three-dimensional variables (Maher et al., 2025).  
973 Both the MMLEAv2 and a suite of corresponding observational datasets have been regridded onto a 2.5° common horizontal  
974 grid, reducing data volume and enabling straightforward model-to-model or model-to-observation comparisons. In addition,  
975 the release of MMLEAv2 is accompanied by the newest version of the CVDP (CVDPv6; Phillips et al., 2020), introduced in  
976 Subsection 2.5.



977

978 Figure 5. Number of publications per year involving SMILEs. The data was extracted from the citation report available at  
979 Web of Science (<https://www.webofscience.com/wos/woscc/basic-search>) for the queries provided in Appendix A.

980 For a variable with relatively low internal variability or high signal-to-noise ratio, 10 ensemble members can be used to  
981 sufficiently detect changes, e.g. for the global mean land temperature (Deser et al., 2012b). To robustly detect significant  
982 warming in the 2050s relative to the 2010s (at the 95% confidence level), Deser et al. (2012b) found that only a single ensemble  
983 member was required for nearly all locations. In contrast, precipitation exhibits substantially higher internal variability:  
984 detecting changes requires approximately 3–6 ensemble members in the tropics and high latitudes, while more than 15

985 ensemble members are needed in the mid-latitudes, with estimates reaching up to 40 members (Deser et al., 2012a, b). For sea  
 986 level pressure, Deser et al. (2012b) reported that only 3-6 ensemble members are sufficient in the tropics, whereas 9-30 are  
 987 required in the extratropics. The number of ensemble members required varies regionally, reflecting differences in both the  
 988 strength of the forced signal and local internal variability (Bittner et al., 2016). Over the ocean, less SMILE members are  
 989 required (Milinski et al., 2020). Table 2 provides a small sample of papers that have employed large ensembles for a variety  
 990 of research questions.

991 **Table 2.** Examples of large ensembles used and how many models were investigated.

Variable/Metric	No. of ensemble members	Study
Aridity and risk of consecutive drought years	Two 10-member ensembles from CESM	(Lehner et al., 2017)
Precipitation and temperature	Two 10-member atmosphere only ensembles from CESM and GFDL 40 models (1 simulation each) from CMIP5 40-member CESM1 Large Ensemble 10-member GFDL Large Ensemble	(Lehner et al., 2018)
Ocean carbon uptake	38-member CESM1-LE 9 models from CMIP5	(Lovenduski et al., 2016)
Temperature and precipitation influence on snow trends	40-member CESM1-LE	(Mankin and Diffenbaugh, 2015)
Irreducible uncertainty	100-member MPI Grand Ensemble	(Marotzke, 2019)
Ocean ecosystem drivers	30-member GFDL Ensemble	(Rodgers et al., 2015)
Ocean carbon cycle	30-member GFDL Ensemble	(Schlunegger et al., 2019)
Weather regimes and their impact on surface extremes	100-member ensemble from CESM2-LE and 36-member ensemble from E3SM2-LE	(Pérez-Carrasquilla et al., 2025)

992

### 993 4.3 Computational Resources and Carbon Impact

994 MMEs, such as CMIP6, are powerful tools for exploring past, current, and future climate conditions, but they come with  
 995 substantial computational and energy demands. MMEs typically rely on multiple simulations across different models or

996 multiple ensemble members of a single model, performed on high-performance computing (HPC) platforms. These execute  
 997 calculations across many parallel cores and process and generate large amounts of data that require careful management and  
 998 optimization. Simulating a century-scale global climate model with high spatial and temporal resolutions can take weeks, even  
 999 on HPC systems. For example, the MPI-ESM1.2 model in its standard low-resolution configuration (approximately 200 km  
 1000 grid spacing), achieves between roughly 45 and 85 simulated years per day, representing a significant improvement over the  
 1001 17 years per day achieved during CMIP5 simulations (Mauritsen et al., 2019). On the other hand, running an ultrahigh-  
 1002 resolution climate model in a near-global setup, with ~1 km horizontal resolution reaches a performance of only about 0.043  
 1003 simulated years per day (~15.7 simulated days per day) (Fuhrer et al., 2018). Computational performance therefore represents  
 1004 a major constraint in the design of ESM experiments, requiring trade-offs between resolution, complexity, and the size of  
 1005 ensembles.

### 1006 **CPMIP metrics for Climate Modelling**

1007 Balaji et al. (2017) introduced a universal set of metrics to evaluate HPC and ESM performance, emphasising that traditional  
 1008 metrics (e.g., floating point operations per second) are no longer sufficient to characterize newer generations of computing  
 1009 architectures and the diverse structures of modern ESMs. Given the increasing complexity of ESM components and their  
 1010 heterogeneous computational characteristics, they advocated adopting these metrics (Table 3) as a standard within globally  
 1011 coordinated modelling initiatives and proposed their collection through the Computational Performance MIP (CPMIP). These  
 1012 metrics provide a consistent basis for assessing technological advances in climate models, are accessible from routine  
 1013 production runs, and capture efficiency across the entire modelling lifecycle.

1014 **Table 3.** List of metrics introduced in CPMIP, adapted from Acosta et al. (2024).

<b>Metric</b>	<b>Short description of the metric</b>
Resolution (spatial degrees of freedom)	Number of grid points per model component
Complexity	Number of prognostic variables per component
Platform	Description of the computational hardware (core count, clock speed, and double-precision operations per clock cycle)
Simulation years per day (SYPD)	Number of simulated years per day in a 24-hour period on a given platform
Actual SYPD (ASYPD)	Actual simulated years per day for a long-running simulation on a given platform (system interruptions, queue wait time, or issues with the model workflow accounted)
Core hours per simulated year (CHSY)	Cost, measured in core hours per simulated year

Parallelization	Total number of cores allocated for the run
Joules per simulated year (JPSY)	Energy cost per simulated year
Coupling cost	Computing cost of the coupling algorithm and load imbalance
Memory bloat	Ratio of actual memory size to ideal memory size
Data output cost	Computing cost for performing input/output (I/O)
Data intensity	Measure of data produced per computing hour

1015

1016 Addressing the performance characteristics of individual models within a MME can contribute to a more balanced and efficient  
1017 use of computational resources. Building on the foundational CPMIP framework, Acosta et al. (2024) extended this work by  
1018 incorporating empirical data from CMIP6, collected during long, real-time model runs across 14 institutions, encompassing  
1019 33 experiments and nearly 500,000 years of simulations. The study places particular emphasis on energy consumption, data  
1020 storage demands, and operational efficiency, and provides strategic recommendations for improving the sustainability and  
1021 performance of future climate modelling efforts.

1022 The demand for computing power is steadily increasing due to several factors, including higher spatial and temporal resolution,  
1023 the explicit representation of complex climate processes that were previously parameterized, increasing ensemble sizes, and  
1024 the associated growth in data storage requirements for both input and output. Improving model accuracy through these  
1025 processes leads to more detailed spatial and temporal outputs, but requires immense computational resources. For example,  
1026 Flato (2011) found that increasing model resolution from 200 km to 20 km demands roughly 10,000 times more computing  
1027 power.

1028 Kilometer-scale simulations of individual models and associated MMEs are being actively developed (Ban et al., 2021;  
1029 Coppola et al., 2020; Pichelli et al., 2021; Rackow et al., 2025), as well as coordinated intercomparisons for global storm-  
1030 resolving models (GSRM). To cope with the high computational and energy demands required for increase in resolution and  
1031 process detail (Schär et al., 2020), such simulations are often conducted over limited regional domains (Coppola et al., 2020;  
1032 Nolan and Flanagan, 2020), rely on simplified parameterizations (for processes such as radiation or soil interactions) or—  
1033 when performed globally—are restricted to relatively short simulation periods of only a few weeks (Schär et al., 2020).  
1034 However, this limitation is rapidly being overcome, with multi-year global simulations at such resolutions have already been  
1035 conducted using models such as ICON in its Sapphire configuration (Hohenegger et al., 2023), the eXperimental System for  
1036 High-resolution prediction on Earth-to-Local Domains (X-SHiELD) (Guendelman et al., 2024; Merlis et al., 2024), or the IFS  
1037 model coupled to the Finite-volume Sea ice-Ocean Model (Rackow et al., 2025).

1038 Beyond the role of resolution and resolved processes, the internal structure of ESMs also plays a crucial role in determining  
1039 computational performance. ESMs are structured with a component-based architecture, allowing scientists to update or add  
1040 new components over time. While this architecture enables continuous flexibility, it also brings software engineering  
1041 challenges. Modifications to individual components can alter computational demands and affect data processing, input/output  
1042 (I/O) operations, and network traffic (Wang and Yuan, 2020). As shown in Acosta et al. (2024) coupling components, which  
1043 synchronize different processes, adds up to 5–15% overhead to execution costs.

1044 Operational factors such as job scheduling also influence computational efficiency. Queue times significantly impact overall  
1045 execution speed and efficiency, although they vary across different institutions (Acosta et al., 2024). Short and consistent  
1046 queue times are beneficial for MMEs, as they help ensure the timely completion of simulations. Reducing queue times  
1047 improves the effective use of available HPC resources, allowing more simulations to be completed within a given timeframe  
1048 and thereby increasing the overall throughput of the ensemble.

### 1049 **Carbon Footprint of Climate Modeling: Towards "greener" Hardware**

1050 Running climate models on HPCs, particularly for large-scale MMEs, requires substantial energy and is associated with a  
1051 significant carbon footprint. The climate modelling community is aware of this and is exploring ways to optimize code  
1052 efficiency and transition to greener energy sources to minimise the carbon impact of their research efforts. In this context,  
1053 CPMIP focuses on capturing the real energy costs of running models, aiming to help climate scientists make eco-friendly  
1054 decisions in computing. Using the CPMIP metrics and the efforts of the Infrastructure for the European Network for Earth  
1055 System Modelling Phase 3 (IS-ENES3) project (Joussaume and Budich, 2013), the total computational energy costs of climate  
1056 experiments have been assessed, enabling estimates of the associated carbon footprint (Acosta et al., 2024). For 8 out of 49  
1057 institutions that were involved in CMIP6, the total estimated emissions amount to 1,692 t CO<sub>2</sub> (with total energy costs ranging  
1058 from 0.41 TJ to 26.70 TJ). For context, the International Energy Agency (IEA) reports that the “global average energy-related  
1059 carbon footprint” is approximately 4.7 t CO<sub>2</sub> per person and per year. Thus, the total emissions from this share of CMIP6  
1060 modelling centers are roughly equivalent to the annual energy-related emissions of 360 people.

1061 As researchers recognize the environmental impact of extensive model runs, eco-friendly hardware is becoming an increasingly  
1062 important consideration in HPC for climate modelling. One example of this good practice is the Energy-efficient climate  
1063 simulations on heterogeneous supercomputers through co-design (EEcliPs) project led by German Climate Computing Centre  
1064 (Deutsches Klimarechenzentrum, DKRZ; <https://www.dkrz.de/en/projects-and-partners/projects-1/eeclips>), aiming to  
1065 improve simulation quality while reducing the energy requirements of the ESM ICON (Adamidis et al., 2025). By encouraging  
1066 institutions to collect the data needed to estimate their carbon footprint, to adopt eco-friendly hardware and to implement  
1067 thoughtful modelling practices, the climate modelling community can reduce its carbon impact while continuing to advance  
1068 its scientific mission.

## 1069 **HPC Facilities: petascale and beyond**

1070 As climate models continue to evolve, HPC facilities operating at the petascale and beyond are necessary. Looking ahead,  
1071 exascale computing systems, capable of achieving  $10^{18}$  floating-point operations per second, promise to greatly expand  
1072 modelling capabilities, enabling longer, higher-resolution simulations, more complex process representation, and improved  
1073 exploration of predictability limits in ESMs. This potential led to the launch of many projects aiming to develop and optimize  
1074 the parallel execution on exascale systems (Adamidis et al., 2025; Taylor et al., 2023; [https://www.fz-  
1075 juelich.de/en/ias/jsc/projects/ifces2](https://www.fz-juelich.de/en/ias/jsc/projects/ifces2)).

1076 In the context of the increasing computational complexity, the findings from the CPMIP and performance metrics applied to  
1077 CMIP6 experiments underline the need for better optimization of model configurations, improved coupling mechanisms, and  
1078 more efficient use of HPC resources. The intercomparison between models and institutions reveals significant differences in  
1079 computational costs, highlighting the need and potential for strategic advancements. Coordinated efforts are also required to  
1080 integrate the latest technological advances, as e.g. outlined in Subsections 4.1 and 4.2. Standardized measurements of  
1081 computational and energy costs can guide the path forward, ensuring that model performance is comparable, and thereby  
1082 allowing researchers to identify areas for improvement and make informed decisions in the development.

## 1083 **5. Concluding Remarks**

1084 Climate modelling has been key to the understanding of past, present, and future climate change. It is a dynamic field, profiting  
1085 from growing computational capacities and advances as well as benefits from the increasing understanding of physical and  
1086 chemical phenomena. Climate projections rely on MMEs to assess uncertainties and improve their robustness. This review  
1087 synthesizes key practices, challenges, and emerging approaches in working with MMEs, drawing on the collective insights of  
1088 the Fresh Eyes on CMIP community. By examining model evaluation strategies, model dependence, selection and weighting  
1089 methods, and uncertainty quantification, we aim to support researchers in making informed choices when designing MME  
1090 studies—while fully acknowledging that the diversity of research questions makes it impossible to create a set of universally  
1091 transferable recommendations. We further highlight the growing relevance of ML and SMILEs, which are shaping the future  
1092 of climate ensemble analysis, particularly in the context of CMIP7. Finally, we advocate for awareness of the computational  
1093 costs associated with climate modelling and analyses.

## 1094 **Acknowledgements**

1095 We thank the wider Fresh Eyes on CMIP community and steering group, the CMIP International Project Office as well as the broader CMIP  
1096 community for their valuable engagement and support, as well as for providing foundational infrastructure including communication  
1097 platforms that made this work possible. We specifically thank Elisabeth Dingley and Yuhan Douglas Rao for their valuable guidance and  
1098 continuous support. We greatly acknowledge the valuable feedback provided by Ranjini Swaminathan and Tomoki Miyakawa during the  
1099 CMIP internal review process, which helped improve the manuscript. We also want to specifically thank the two anonymous reviewers that  
1100 contributed to a substantially improved version of this manuscript. We acknowledge Josh Dorrington for initial help with this project. NČ  
1101 thanks Robert Pincus, Gregory Cesana, Andrew Ackerman, and others at Columbia CCSR and NASA GISS for introducing her to the CMIP  
1102 community and for insightful discussions on various topics related to analysis and evaluation of CMIP MMEs in earlier projects. JSPC  
1103 thanks Maria J. Molina for mentoring on how ML can assist the climate science community, and Isla R. Simpson for sharing her expertise  
1104 in uncertainty in climate projections and emergent constraints. AK thanks Anders Levermann, Jacob Schewe, and Julia Pongratz for  
1105 insightful discussions on MME design in earlier projects, which offered a valuable foundation for the present study. MT thanks Vladimir  
1106 Djurdjevic for his foundational mentorship in understanding global and regional climate model ensembles, and Theodore Shepherd for  
1107 insightful discussions that significantly shaped her thinking on uncertainty, storylines and the interpretation of MME data. EG thanks  
1108 Veronika Eyring and Jakob Runge for useful discussions on ML and causality topics related to climate model evaluation. CL thanks Kirsten  
1109 Zickfeld for valuable exchanges on result robustness and statistical analysis more generally and Alex Koch for the introduction to working  
1110 with CMIP data.

1111 CL acknowledges support from the Natural Sciences and Engineering Research Council of Canada Discovery Grant Program (grant no.  
1112 RGPIN-2018-06881 awarded to K. Zickfeld. EG is funded by the Central Research Development Fund at the University of Bremen, Funding  
1113 No: ZF04A/2023/FB1/Galytska Evgenia. PP acknowledges the financial support received in the form of a doctoral research fellowship from  
1114 the Council of Scientific and Industrial Research (CSIR), India, Award no: 09/1187(11135)/2021-EMR-I. M. T. acknowledges support from  
1115 the Science Fund of the Republic of Serbia (Grant No. 7389, Project Extreme weather events in Serbia - analysis, modelling and impacts” -  
1116 EXTREMES). JSPC was supported by a University of Maryland Grand Challenges Seed Grant. NČ acknowledges support from the NOAA  
1117 grant NA20OAR4310390, the NASA Modeling, Analysis, and Prediction Program number 80NSSC21K1134, ARIS Programme P1-0188,  
1118 and the University of Ljubljana Grant SN-ZRD/22-27/0510, which covers the fee costs of this publication.

## 1119 **Author Contribution Statement**

1120 All authors conducted a literature review, contributed valuable ideas to the scientific content and study design, topic discussions, and writing  
1121 of the manuscript (Abstract: AK, NČ; Introduction: NČ, AK; Subsection 2.1: EG, AK, IR, NČ; Subsection 2.2: KG; Subsection 2.3: KG,  
1122 PP; Subsection 2.4: JSPC, MT; Subsection 2.5: NČ, EG, MT; Subsection 3.1: AK; Subsection 3.2: MT; Subsection 3.3: CL; Subsection 3.4:  
1123 PP; Subsection 3.5: CL; Subsection 4.1: EG, KG, JSPC; Section 4.2: AVC, MT, AK; Subsection 4.3: MT; Conclusion: AK; Appendix B:  
1124 NČ; Appendix D: JSPC). Final details will be provided with publication.

1126 **References**

- 1127 Abramowitz, G., Herger, N., Gutmann, E., Hammerling, D., Knutti, R., Leduc, M., Lorenz, R., Pincus, R., and Schmidt, G.  
 1128 A.: ESD Reviews: Model dependence in multi-model climate ensembles: weighting, sub-selection and out-of-sample testing,  
 1129 *Earth Syst. Dyn.*, 10, 91–105, <https://doi.org/10.5194/esd-10-91-2019>, 2019.
- 1130 Achugbu, I. C., Olufayo, A. A., Balogun, I. A., Adefisan, E. A., Dudhia, J., and Naabil, E.: Modeling the spatiotemporal  
 1131 response of dew point temperature, air temperature and rainfall to land use land cover change over West Africa, *Model.*  
 1132 *Earth Syst. Environ.*, 8, 173–198, <https://doi.org/10.1007/s40808-021-01094-8>, 2022.
- 1133 Acosta, M. C., Palomas, S., Paronuzzi Ticco, S. V., Utrera, G., Biercamp, J., Bretonniere, P.-A., Budich, R., Castrillo, M.,  
 1134 Caubel, A., Doblas-Reyes, F., Epicoco, I., Fladrich, U., Joussaume, S., Kumar Gupta, A., Lawrence, B., Le Sager, P., Lister,  
 1135 G., Moine, M.-P., Rioual, J.-C., Valcke, S., Zadeh, N., and Balaji, V.: The computational and energy cost of simulation and  
 1136 storage for climate science: lessons from CMIP6, *Geosci. Model Dev.*, 17, 3081–3098, <https://doi.org/10.5194/gmd-17-3081-2024>, 2024.
- 1138 Adamidis, P., Pfister, E., Bockelmann, H., Zobel, D., Beismann, J.-O., and Jacob, M.: The real challenges for climate and  
 1139 weather modelling on its way to sustained exascale performance: a case study using ICON (v2.6.6), *Geosci. Model Dev.*, 18,  
 1140 905–919, <https://doi.org/10.5194/gmd-18-905-2025>, 2025.
- 1141 Ahn, M., Daehyun, K., Sperber, K. R., Kang, I.-S., Maloney, E., Waliser, D., Hendon, H., and on behalf of WGNE MJO  
 1142 Task Force: MJO simulation in CMIP5 climate models: MJO skill metrics and process-oriented diagnosis, *Clim. Dyn.*, 49,  
 1143 4023–4045, <https://doi.org/10.1007/s00382-017-3558-4>, 2017.
- 1144 Ahn, M., Kim, D., Kang, D., Lee, J., Sperber, K. R., Gleckler, P. J., Jiang, X., Ham, Y., and Kim, H.: MJO Propagation  
 1145 Across the Maritime Continent: Are CMIP6 Models Better Than CMIP5 Models?, *Geophys. Res. Lett.*, 47,  
 1146 e2020GL087250, <https://doi.org/10.1029/2020GL087250>, 2020.
- 1147 Almazroui, M., Saeed, S., Islam, M. N., Khalid, M. S., Alkhalaf, A. K., and Dambul, R.: Assessment of uncertainties in  
 1148 projected temperature and precipitation over the Arabian Peninsula: a comparison between different categories of CMIP3  
 1149 models, *Earth Syst. Environ.*, 1, 12, <https://doi.org/10.1007/s41748-017-0012-z>, 2017.
- 1150 Amali, A. A., Schwingshackl, C., Ito, A., Barbu, A., Delire, C., Peano, D., Lawrence, D. M., Wårlind, D., Robertson, E.,  
 1151 Davin, E. L., Shevliakova, E., Harman, I. N., Vuichard, N., Miller, P. A., Lawrence, P. J., Ziehn, T., Hajima, T., Brovkin, V.,  
 1152 Zhang, Y., Arora, V. K., and Pongratz, J.: Biogeochemical versus biogeophysical temperature effects of historical land-use  
 1153 change in CMIP6, <https://doi.org/10.5194/egusphere-2024-2460>, 27 August 2024.
- 1154 Annan, J. D. and Hargreaves, J. C.: On the meaning of independence in climate science, *Earth Syst. Dyn.*, 8, 211–224,  
 1155 <https://doi.org/10.5194/esd-8-211-2017>, 2017.
- 1156 NetCDF Users Guide: NetCDF Utilities: [https://docs.unidata.ucar.edu/nug/current/netcdf\\_utilities\\_guide.html](https://docs.unidata.ucar.edu/nug/current/netcdf_utilities_guide.html), last access:  
 1157 12 May 2025.
- 1158 Aru, H., Chen, W., Chen, S., Garfinkel, C. I., Ma, T., Dong, Z., and Hu, P.: Variation in the Impact of ENSO on the Western  
 1159 Pacific Pattern Influenced by ENSO Amplitude in CMIP6 Simulations, *J. Geophys. Res. Atmospheres*, 128,  
 1160 e2022JD037905, <https://doi.org/10.1029/2022JD037905>, 2023.
- 1161 Balaji, V., Maisonnave, E., Zadeh, N., Lawrence, B. N., Biercamp, J., Fladrich, U., Aloisio, G., Benson, R., Caubel, A.,  
 1162 Durachta, J., Foujols, M.-A., Lister, G., Mocavero, S., Underwood, S., and Wright, G.: CPMIP: measurements of real  
 1163 computational performance of Earth system models in CMIP6, *Geosci. Model Dev.*, 10, 19–34, <https://doi.org/10.5194/gmd-10-19-2017>, 2017.
- 1165 Balhane, S., Driouech, F., Chafki, O., Manzanar, R., Chehbouni, A., and Moufouma-Okia, W.: Changes in mean and  
 1166 extreme temperature and precipitation events from different weighted multi-model ensembles over the northern half of

- 1167 Morocco, *Clim. Dyn.*, 58, 389–404, <https://doi.org/10.1007/s00382-021-05910-w>, 2022.
- 1168 Ban, N., Caillaud, C., Coppola, E., Pichelli, E., Sobolowski, S., Adinolfi, M., Ahrens, B., Alias, A., Anders, I., Bastin, S.,  
1169 Belušić, D., Berthou, S., Brisson, E., Cardoso, R. M., Chan, S. C., Christensen, O. B., Fernández, J., Fita, L., Frisius, T.,  
1170 Gašparac, G., Giorgi, F., Goergen, K., Haugen, J. E., Hodnebrog, Ø., Kartsios, S., Katragkou, E., Kendon, E. J., Keuler, K.,  
1171 Lavin-Gullon, A., Lenderink, G., Leutwyler, D., Lorenz, T., Maraun, D., Mercogliano, P., Milovac, J., Panitz, H.-J., Raffa,  
1172 M., Remedio, A. R., Schär, C., Soares, P. M. M., Srnec, L., Steensen, B. M., Stocchi, P., Tölle, M. H., Truhetz, H., Vergara-  
1173 Temprado, J., de Vries, H., Warrach-Sagi, K., Wulfmeyer, V., and Zander, M. J.: The first multi-model ensemble of regional  
1174 climate simulations at kilometer-scale resolution, part I: evaluation of precipitation, *Clim. Dyn.*, 57, 275–302,  
1175 <https://doi.org/10.1007/s00382-021-05708-w>, 2021.
- 1176 Becker, E., Kirtman, B. P., and Pegion, K.: Evolution of the North American Multi-Model Ensemble, *Geophys. Res. Lett.*,  
1177 47, e2020GL087408, <https://doi.org/10.1029/2020GL087408>, 2020.
- 1178 Becker, E. J., Kirtman, B. P., L’Heureux, M., Muñoz, Á. G., and Pegion, K.: A Decade of the North American Multimodel  
1179 Ensemble (NMME): Research, Application, and Future Directions, *Bull. Am. Meteorol. Soc.*, 103, E973–E995,  
1180 <https://doi.org/10.1175/BAMS-D-20-0327.1>, 2022.
- 1181 Bellomo, K., Angeloni, M., Corti, S., and von Hardenberg, J.: Future climate change shaped by inter-model differences in  
1182 Atlantic meridional overturning circulation response, *Nat. Commun.*, 12, 3659, [https://doi.org/10.1038/s41467-021-24015-](https://doi.org/10.1038/s41467-021-24015-w)  
1183 w, 2021.
- 1184 Beucler, T., Gentine, P., Yuval, J., Gupta, A., Peng, L., Lin, J., Yu, S., Rasp, S., Ahmed, F., O’Gorman, P. A., Neelin, J. D.,  
1185 Lutsko, N. J., and Pritchard, M.: Climate-invariant machine learning, *Sci. Adv.*, 10, eadj7250,  
1186 <https://doi.org/10.1126/sciadv.adj7250>, 2024.
- 1187 Bevacqua, E., Zappa, G., Lehner, F., and Zscheischler, J.: Precipitation trends determine future occurrences of compound  
1188 hot–dry events, *Nat. Clim. Change*, 12, 350–355, <https://doi.org/10.1038/s41558-022-01309-5>, 2022.
- 1189 Bevacqua, E., Suarez-Gutierrez, L., Jézéquel, A., Lehner, F., Vrac, M., Yiou, P., and Zscheischler, J.: Advancing research on  
1190 compound weather and climate events via large ensemble model simulations, *Nat Commun*, 14, 2145,  
1191 <https://doi.org/10.1038/s41467-023-37847-5>, 2023.
- 1192 Beven, K. and Freer, J.: Equifinality, data assimilation, and uncertainty estimation in mechanistic modelling of complex  
1193 environmental systems using the GLUE methodology, *J. Hydrol.*, 249, 11–29, [https://doi.org/10.1016/S0022-](https://doi.org/10.1016/S0022-1694(01)00421-8)  
1194 1694(01)00421-8, 2001.
- 1195 Bhowmik, R. and Sankarasubramanian, A.: A performance-based multi-model combination approach to reduce uncertainty  
1196 in seasonal temperature change projections, *Int. J. Climatol.*, 41, <https://doi.org/10.1002/joc.6870>, 2020.
- 1197 Bittner, M., Schmidt, H., Timmreck, C., and Sienz, F.: Using a large ensemble of simulations to assess the Northern  
1198 Hemisphere stratospheric dynamical response to tropical volcanic eruptions and its uncertainty, *Geophys. Res. Lett.*, 43,  
1199 9324–9332, <https://doi.org/10.1002/2016GL070587>, 2016.
- 1200 Boé, J.: Interdependency in Multimodel Climate Projections: Component Replication and Result Similarity, *Geophys. Res.*  
1201 *Lett.*, 45, 2771–2779, <https://doi.org/10.1002/2017GL076829>, 2018.
- 1202 Boysen, L. R.: BG - Global climate response to idealized deforestation in CMIP6 models, 2020.
- 1203 Bracegirdle, T. J. and Stephenson, D. B.: Higher precision estimates of regional polar warming by ensemble regression of  
1204 climate model projections, *Clim. Dyn.*, 39, 2805–2821, <https://doi.org/10.1007/s00382-012-1330-3>, 2012.
- 1205 Brenowitz, N. D., Henn, B., McGibbon, J., Clark, S. K., Kwa, A., Perkins, W. A., Watt-Meyer, O., and Bretherton, C. S.:  
1206 Machine Learning Climate Model Dynamics: Offline versus Online Performance,  
1207 <https://doi.org/10.48550/ARXIV.2011.03081>, 2020.
- 1208 Breul, P., Ceppi, P., and Shepherd, T. G.: Revisiting the wintertime emergent constraint of the southern hemispheric  
1209 midlatitude jet response to global warming, *Weather Clim. Dyn.*, 4, 39–47, <https://doi.org/10.5194/wcd-4-39-2023>, 2023.

- 1210 Brunner, L. and Sippel, S.: Identifying climate models based on their daily output using machine learning, *Environ. Data*  
1211 *Sci.*, 2, e22, <https://doi.org/10.1017/eds.2023.23>, 2023.
- 1212 Brunner, L., Pendergrass, A. G., Lehner, F., Merrifield, A. L., Lorenz, R., and Knutti, R.: Reduced global warming from  
1213 CMIP6 projections when weighting models by performance and independence, *Earth Syst. Dyn.*, 11, 995–1012,  
1214 <https://doi.org/10.5194/esd-11-995-2020>, 2020.
- 1215 Buontempo, C., Burgess, S. N., Dee, D., Pinty, B., Thépaut, J.-N., Rixen, M., Almond, S., Armstrong, D., Brookshaw, A.,  
1216 Alos, A. L., Bell, B., Bergeron, C., Cagnazzo, C., Comyn-Platt, E., Damasio-Da-Costa, E., Guillory, A., Hersbach, H.,  
1217 Horányi, A., Nicolas, J., Obregon, A., Ramos, E. P., Raoult, B., Muñoz-Sabater, J., Simmons, A., Soci, C., Suttie, M.,  
1218 Vamborg, F., Varndell, J., Vermoote, S., Yang, X., and Garcés De Marcilla, J.: The Copernicus Climate Change Service:  
1219 Climate Science in Action, *Bull. Am. Meteorol. Soc.*, 103, E2669–E2687, <https://doi.org/10.1175/BAMS-D-21-0315.1>,  
1220 2022.
- 1221 de Burgh-Day, C. O. and Leeuwenburg, T.: Machine learning for numerical weather and climate modelling: a review,  
1222 *Geosci. Model Dev.*, 16, 6433–6477, <https://doi.org/10.5194/gmd-16-6433-2023>, 2023.
- 1223 Cesana, G. V. and Del Genio, A. D.: Observational constraint on cloud feedbacks suggests moderate climate sensitivity, *Nat.*  
1224 *Clim. Change*, 11, 213–218, <https://doi.org/10.1038/s41558-020-00970-y>, 2021.
- 1225 Cesana, G. V., Ackerman, A. S., Črnivec, N., Pincus, R., and Chepfer, H.: An observation-based method to assess tropical  
1226 stratocumulus and shallow cumulus clouds and feedbacks in CMIP6 and CMIP5 models, *Environ. Res. Commun.*, 5,  
1227 045001, <https://doi.org/10.1088/2515-7620/acc78a>, 2023.
- 1228 Chandra, S., Kumar, P., Siingh, D., Roy, I., Victor, N. J., and Kamra, A. K.: Projection of lightning over South/South East  
1229 Asia using CMIP5 models, *Nat. Hazards*, 114, 57–75, <https://doi.org/10.1007/s11069-022-05379-8>, 2022.
- 1230 Cinquini, L., Crichton, D., Mattmann, C., Harney, J., Shipman, G., Wang, F., Ananthakrishnan, R., Miller, N., Denvil, S.,  
1231 Morgan, M., Pobre, Z., Bell, G. M., Drach, B., Williams, D., Kershaw, P., Pascoe, S., Gonzalez, E., Fiore, S., and  
1232 Schweitzer, R.: The Earth System Grid Federation: An open infrastructure for access to distributed geospatial data, in: 2012  
1233 IEEE 8th International Conference on E-Science, 2012 IEEE 8th International Conference on E-Science, 1–10,  
1234 <https://doi.org/10.1109/eScience.2012.6404471>, 2012.
- 1235 Clyde, M., Çetinkaya-Rundel, M., Rundel, C., Banks, D., Chai, C., and Huang, L.: An Introduction to Bayesian Thinking,  
1236 2022.
- 1237 Coles, S.: *An Introduction to Statistical Modeling of Extreme Values*, Springer London, London,  
1238 <https://doi.org/10.1007/978-1-4471-3675-0>, 2001.
- 1239 Cook, B. I., Mankin, J. S., Marvel, K., Williams, A. P., Smerdon, J. E., and Anchukaitis, K. J.: Twenty-First Century  
1240 Drought Projections in the CMIP6 Forcing Scenarios, *Earths Future*, 8, e2019EF001461,  
1241 <https://doi.org/10.1029/2019EF001461>, 2020.
- 1242 Coppola, E., Sobolowski, S., Pichelli, E., Raffaele, F., Ahrens, B., Anders, I., Ban, N., Bastin, S., Belda, M., Belusic, D.,  
1243 Caldas-Alvarez, A., Cardoso, R. M., Davolio, S., Dobler, A., Fernandez, J., Fita, L., Fumiere, Q., Giorgi, F., Goergen, K.,  
1244 Güttler, I., Halenka, T., Heinzeller, D., Hodnebrog, Ø., Jacob, D., Kartsios, S., Katragkou, E., Kendon, E., Khodayar, S.,  
1245 Kunstmann, H., Knist, S., Lavín-Gullón, A., Lind, P., Lorenz, T., Maraun, D., Marelle, L., van Meijgaard, E., Milovac, J.,  
1246 Myhre, G., Panitz, H.-J., Piazza, M., Raffa, M., Raub, T., Rockel, B., Schär, C., Sieck, K., Soares, P. M. M., Somot, S.,  
1247 Srncic, L., Stocchi, P., Tölle, M. H., Truhetz, H., Vautard, R., de Vries, H., and Warrach-Sagi, K.: A first-of-its-kind multi-  
1248 model convection permitting ensemble for investigating convective phenomena over Europe and the Mediterranean, *Clim.*  
1249 *Dyn.*, 55, 3–34, <https://doi.org/10.1007/s00382-018-4521-8>, 2020.
- 1250 Coppola, E., Nogherotto, R., Ciarlo, J. M., Giorgi, F., Van Meijgaard, E., Kadyrov, N., Iles, C., Corre, L., Sandstad, M.,  
1251 Somot, S., Nabat, P., Vautard, R., Levavasseur, G., Schwingshackl, C., Sillmann, J., Kjellström, E., Nikulin, G., Aalbers, E.,  
1252 Lenderink, G., Christensen, O. B., Boberg, F., Sørland, S. L., Demory, M., Bülow, K., Teichmann, C., Warrach-Sagi, K., and  
1253 Wulfmeyer, V.: Assessment of the European Climate Projections as Simulated by the Large EURO-CORDEX Regional and  
1254 Global Climate Model Ensemble, *J. Geophys. Res. Atmospheres*, 126, e2019JD032356,

- 1255 <https://doi.org/10.1029/2019JD032356>, 2021.
- 1256 Crane-Droesch, A.: Machine learning methods for crop yield prediction and climate change impact assessment in  
1257 agriculture, *Environ. Res. Lett.*, 13, 114003, <https://doi.org/10.1088/1748-9326/aae159>, 2018.
- 1258 Crawford, J., Venkataraman, K., and Booth, J.: Developing climate model ensembles: A comparative case study, *J. Hydrol.*,  
1259 568, 160–173, <https://doi.org/10.1016/j.jhydrol.2018.10.054>, 2019.
- 1260 Črnivec, N., Cesana, G., and Pincus, R.: Evaluating the Representation of Tropical Stratocumulus and Shallow Cumulus  
1261 Clouds As Well As Their Radiative Effects in CMIP6 Models Using Satellite Observations, *J. Geophys. Res. Atmospheres*,  
1262 128, e2022JD038437, <https://doi.org/10.1029/2022JD038437>, 2023.
- 1263 Debeire, K., Bock, L., Nowack, P., Runge, J., and Eyring, V.: Constraining uncertainty in projected precipitation over land  
1264 with causal discovery, *Earth Syst. Dyn.*, 16, 607–630, <https://doi.org/10.5194/esd-16-607-2025>, 2025.
- 1265 DelSole, T. and Tippet, M.: *Statistical Methods for Climate Scientists*, Cambridge University Press, Cambridge,  
1266 <https://doi.org/10.1017/9781108659055>, 2022.
- 1267 Deser, C.: “Certain Uncertainty: The Role of Internal Climate Variability in Projections of Regional Climate Change and  
1268 Risk Management,” *Earths Future*, 8, e2020EF001854, <https://doi.org/10.1029/2020EF001854>, 2020.
- 1269 Deser, C., Knutti, R., Solomon, S., and Phillips, A. S.: Communication of the role of natural variability in future North  
1270 American climate, *Nat. Clim. Change*, 2, 775–779, <https://doi.org/10.1038/nclimate1562>, 2012a.
- 1271 Deser, C., Phillips, A., Bourdette, V., and Teng, H.: Uncertainty in climate change projections: the role of internal  
1272 variability, *Clim. Dyn.*, 38, 527–546, <https://doi.org/10.1007/s00382-010-0977-x>, 2012b.
- 1273 Deser, C., Lehner, F., Rodgers, K. B., Ault, T., Delworth, T. L., DiNezio, P. N., Fiore, A., Frankignoul, C., Fyfe, J. C.,  
1274 Horton, D. E., Kay, J. E., Knutti, R., Lovenduski, N. S., Marotzke, J., McKinnon, K. A., Minobe, S., Randerson, J., Screen,  
1275 J. A., Simpson, I. R., and Ting, M.: Insights from Earth system model initial-condition large ensembles and future prospects,  
1276 *Nat. Clim. Change*, 10, 277–286, <https://doi.org/10.1038/s41558-020-0731-2>, 2020.
- 1277 Dey, A., Sahoo, D. P., Kumar, R., and Remesan, R.: A multimodel ensemble machine learning approach for CMIP6 climate  
1278 model projections in an Indian River basin, *Int. J. Climatol.*, 42, 9215–9236, <https://doi.org/10.1002/joc.7813>, 2022.
- 1279 Di Luca, A., De Elía, R., and Laprise, R.: Challenges in the Quest for Added Value of Regional Climate Dynamical  
1280 Downscaling, *Curr. Clim. Change Rep.*, 1, 10–21, <https://doi.org/10.1007/s40641-015-0003-9>, 2015.
- 1281 Di Luca, A., De Elía, R., Bador, M., and Argüeso, D.: Contribution of mean climate to hot temperature extremes for present  
1282 and future climates, *Weather Clim. Extrem.*, 28, 100255, <https://doi.org/10.1016/j.wace.2020.100255>, 2020a.
- 1283 Di Luca, A., Pitman, A. J., and de Elía, R.: Decomposing Temperature Extremes Errors in CMIP5 and CMIP6 Models,  
1284 *Geophys. Res. Lett.*, 47, e2020GL088031, <https://doi.org/10.1029/2020GL088031>, 2020b.
- 1285 Di Virgilio, G., Ji, F., Tam, E., Nishant, N., Evans, J. P., Thomas, C., Riley, M. L., Beyer, K., Grose, M. R., Narsey, S., and  
1286 Delage, F.: Selecting CMIP6 GCMs for CORDEX Dynamical Downscaling: Model Performance, Independence, and  
1287 Climate Change Signals, *Earths Future*, 10, e2021EF002625, <https://doi.org/10.1029/2021EF002625>, 2022.
- 1288 Dirkes, C. A., Wing, A. A., Camargo, S. J., and Kim, D.: Process-Oriented Diagnosis of Tropical Cyclones in Reanalyses  
1289 Using a Moist Static Energy Variance Budget, *J. Clim.*, 36, 5293–5317, <https://doi.org/10.1175/JCLI-D-22-0384.1>, 2023.
- 1290 Doblas-Reyes, F. J., Pavan, V., and Stephenson, D. B.: The skill of multi-model seasonal forecasts of the wintertime North  
1291 Atlantic Oscillation, *Clim. Dyn.*, 21, 501–514, <https://doi.org/10.1007/s00382-003-0350-4>, 2003.
- 1292 Doblas-Reyes, F. J., Hagedorn, R., and Palmer, T. N.: The rationale behind the success of multi-model ensembles in seasonal  
1293 forecasting – II. Calibration and combination, *Tellus Dyn. Meteorol. Oceanogr.*, 57, 234,  
1294 <https://doi.org/10.3402/tellusa.v57i3.14658>, 2005.
- 1295 Docquier, D., Vannitsem, S., Ragone, F., Wyser, K., and Liang, X. S.: Causal Links Between Arctic Sea Ice and Its Potential  
1296 Drivers Based on the Rate of Information Transfer, *Geophys. Res. Lett.*, 49, e2021GL095892,

- 1297 <https://doi.org/10.1029/2021GL095892>, 2022.
- 1298 Docquier, D., Massonnet, F., Ragone, F., Sticker, A., Fichet, T., and Vannitsem, S.: Drivers of summer Arctic sea-ice  
1299 extent in CMIP6 large ensembles revealed by information flow, <https://doi.org/10.21203/rs.3.rs-4434953/v1>, 4 June 2024.
- 1300 Dosio, A.: Projections of climate change indices of temperature and precipitation from an ensemble of bias-adjusted high-  
1301 resolution EURO-CORDEX regional climate models, *J. Geophys. Res. Atmospheres*, 121, 5488–5511,  
1302 <https://doi.org/10.1002/2015JD024411>, 2016.
- 1303 Dosio, A.: Projection of temperature and heat waves for Africa with an ensemble of CORDEX Regional Climate Models,  
1304 *Clim. Dyn.*, 49, 493–519, <https://doi.org/10.1007/s00382-016-3355-5>, 2017.
- 1305 Dueben, P. D. and Bauer, P.: Challenges and design choices for global weather and climate models based on machine  
1306 learning, *Geosci. Model Dev.*, 11, 3999–4009, <https://doi.org/10.5194/gmd-11-3999-2018>, 2018.
- 1307 Eidhammer, T., Gettelman, A., Thayer-Calder, K., Watson-Parris, D., Elsaesser, G., Morrison, H., Van Lier-Walqui, M.,  
1308 Song, C., and McCoy, D.: An extensible perturbed parameter ensemble for the Community Atmosphere Model version 6,  
1309 *Geosci. Model Dev.*, 17, 7835–7853, <https://doi.org/10.5194/gmd-17-7835-2024>, 2024.
- 1310 Eyring, V., Harris, N. R. P., Rex, M., Shepherd, T. G., Fahey, D. W., Amanatidis, G. T., Austin, J., Chipperfield, M. P.,  
1311 Dameris, M., Forster, P. M. D. F., Gettelman, A., Graf, H. F., Nagashima, T., Newman, P. A., Pawson, S., Prather, M. J.,  
1312 Pyle, J. A., Salawitch, R. J., Santer, B. D., and Waugh, D. W.: A Strategy for Process-Oriented Validation of Coupled  
1313 Chemistry–Climate Models, *Bull. Am. Meteorol. Soc.*, 86, 1117–1134, <https://doi.org/10.1175/BAMS-86-8-1117>, 2005.
- 1314 Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., and Taylor, K. E.: Overview of the Coupled  
1315 Model Intercomparison Project Phase 6 (CMIP6) experimental design and organization, *Geosci. Model Dev.*, 9, 1937–1958,  
1316 <https://doi.org/10.5194/gmd-9-1937-2016>, 2016.
- 1317 Eyring, V., Bock, L., Lauer, A., Righi, M., Schlund, M., Andela, B., Arnone, E., Bellprat, O., Brötz, B., Caron, L.-P.,  
1318 Carvalhais, N., Cionni, I., Cortesi, N., Crezee, B., Davin, E. L., Davini, P., Debeire, K., De Mora, L., Deser, C., Docquier,  
1319 D., Earnshaw, P., Ehbrecht, C., Gier, B. K., Gonzalez-Reviriego, N., Goodman, P., Hagemann, S., Hardiman, S., Hassler, B.,  
1320 Hunter, A., Kadow, C., Kindermann, S., Koirala, S., Koldunov, N., Lejeune, Q., Lembo, V., Lovato, T., Lucarini, V.,  
1321 Massonnet, F., Müller, B., Pandde, A., Pérez-Zanón, N., Phillips, A., Predoi, V., Russell, J., Sellar, A., Serva, F., Stacke, T.,  
1322 Swaminathan, R., Torralba, V., Vegas-Regidor, J., Von Hardenberg, J., Weigel, K., and Zimmermann, K.: Earth System  
1323 Model Evaluation Tool (ESMValTool) v2.0 – an extended set of large-scale diagnostics for quasi-operational and  
1324 comprehensive evaluation of Earth system models in CMIP, *Geosci. Model Dev.*, 13, 3383–3438,  
1325 <https://doi.org/10.5194/gmd-13-3383-2020>, 2020.
- 1326 Eyring, V., Mishra, V., Griffith, G. P., Chen, L., Keenan, T., Turetsky, M. R., Brown, S., Jotzo, F., Moore, F. C., and Van  
1327 Der Linden, S.: Reflections and projections on a decade of climate science, *Nat. Clim. Change*, 11, 279–285,  
1328 <https://doi.org/10.1038/s41558-021-01020-x>, 2021.
- 1329 Eyring, V., Collins, W. D., Gentine, P., Barnes, E. A., Barreiro, M., Beucler, T., Bocquet, M., Bretherton, C. S., Christensen,  
1330 H. M., Dagon, K., Gagne, D. J., Hall, D., Hammerling, D., Hoyer, S., Iglesias-Suarez, F., Lopez-Gomez, I., McGraw, M. C.,  
1331 Meehl, G. A., Molina, M. J., Monteleoni, C., Mueller, J., Pritchard, M. S., Rolnick, D., Runge, J., Stier, P., Watt-Meyer, O.,  
1332 Weigel, K., Yu, R., and Zanna, L.: Pushing the frontiers in climate modelling and analysis with machine learning, *Nat. Clim.*  
1333 *Change*, 14, 916–928, <https://doi.org/10.1038/s41558-024-02095-y>, 2024.
- 1334 Falkena, S. K. J. and von der Heydt, A. S.: Subpolar Gyre Variability in CMIP6 Models: Is there a Mechanism for  
1335 Bistability?, <https://doi.org/10.48550/ARXIV.2408.16541>, 2024.
- 1336 Flato, G. M.: Earth system models: an overview, *WIREs Clim. Change*, 2, 783–800, <https://doi.org/10.1002/wcc.148>, 2011.
- 1337 Fuhrer, O., Chadha, T., Hoefler, T., Kwasniewski, G., Lapillonne, X., Leutwyler, D., Lüthi, D., Osuna, C., Schär, C.,  
1338 Schulthess, T. C., and Vogt, H.: Near-global climate simulation at 1 km resolution: establishing a performance baseline on  
1339 4888 GPUs with COSMO 5.0, *Geosci. Model Dev.*, 11, 1665–1681, <https://doi.org/10.5194/gmd-11-1665-2018>, 2018.
- 1340 Galytska, E., Weigel, K., Handorf, D., Jaiser, R., Köhler, R., Runge, J., and Eyring, V.: Evaluating Causal Arctic-

1341 Midlatitude Teleconnections in CMIP6, *J. Geophys. Res. Atmospheres*, 128, e2022JD037978,  
1342 <https://doi.org/10.1029/2022JD037978>, 2023.

1343 Gates, W. L.: AN AMS CONTINUING SERIES: GLOBAL CHANGE--AMIP: The Atmospheric Model Intercomparison  
1344 Project, *Bull. Am. Meteorol. Soc.*, 73, 1962–1970, [https://doi.org/10.1175/1520-  
1345 0477\(1992\)073%253C1962:ATAMIP%253E2.0.CO;2](https://doi.org/10.1175/1520-0477(1992)073%253C1962:ATAMIP%253E2.0.CO;2), 1992.

1346 Ge, F., Zhu, S., Luo, H., Zhi, X., and Wang, H.: Future changes in precipitation extremes over Southeast Asia: insights from  
1347 CMIP6 multi-model ensemble, *Environ. Res. Lett.*, 16, 024013, <https://doi.org/10.1088/1748-9326/abd7ad>, 2021.

1348 Gebrechorkos, S., Leyland, J., Slater, L., Wortmann, M., Ashworth, P. J., Bennett, G. L., Boothroyd, R., Cloke, H., Delorme,  
1349 P., Griffith, H., Hardy, R., Hawker, L., McLelland, S., Neal, J., Nicholas, A., Tatem, A. J., Vahidi, E., Parsons, D. R., and  
1350 Darby, S. E.: A high-resolution daily global dataset of statistically downscaled CMIP6 models for climate impact analyses,  
1351 *Sci. Data*, 10, 611, <https://doi.org/10.1038/s41597-023-02528-x>, 2023.

1352 Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., and Yacalis, G.: Could Machine Learning Break the Convection  
1353 Parameterization Deadlock?, *Geophys. Res. Lett.*, 45, 5742–5751, <https://doi.org/10.1029/2018GL078202>, 2018.

1354 Gergel, D. R., Malevich, S. B., McCusker, K. E., Tenezakis, E., Delgado, M. T., Fish, M. A., and Kopp, R. E.: Global  
1355 Downscaled Projections for Climate Impacts Research (GDPCIR): preserving quantile trends for modeling future climate  
1356 impacts, *Geosci. Model Dev.*, 17, 191–227, <https://doi.org/10.5194/gmd-17-191-2024>, 2024.

1357 Gettelman, A., Geer, A. J., Forbes, R. M., Carmichael, G. R., Feingold, G., Posselt, D. J., Stephens, G. L., Van Den Heever,  
1358 S. C., Varble, A. C., and Zuidema, P.: The future of Earth system prediction: Advances in model-data fusion, *Sci. Adv.*, 8,  
1359 eabn3488, <https://doi.org/10.1126/sciadv.abn3488>, 2022.

1360 Giorgi, F.: Thirty Years of Regional Climate Modeling: Where Are We and Where Are We Going next?, *J. Geophys. Res.*  
1361 *Atmospheres*, 124, 5696–5723, <https://doi.org/10.1029/2018JD030094>, 2019.

1362 Giorgi, F. and Gutowski, W. J.: Regional Dynamical Downscaling and the CORDEX Initiative, *Annu. Rev. Environ.*  
1363 *Resour.*, 40, 467–490, <https://doi.org/10.1146/annurev-environ-102014-021217>, 2015.

1364 Giorgi, F. and Mearns, L. O.: Calculation of Average, Uncertainty Range, and Reliability of Regional Climate Changes from  
1365 AOGCM Simulations via the “Reliability Ensemble Averaging” (REA) Method, *J. Clim.*, 15, 1141–1158,  
1366 [https://doi.org/10.1175/1520-0442\(2002\)015%253C1141:COAURA%253E2.0.CO;2](https://doi.org/10.1175/1520-0442(2002)015%253C1141:COAURA%253E2.0.CO;2), 2002.

1367 Gleckler, P. J., Taylor, K. E., and Doutriaux, C.: Performance metrics for climate models, *J. Geophys. Res. Atmospheres*,  
1368 113, <https://doi.org/10.1029/2007JD008972>, 2008.

1369 Glymour, C., Zhang, K., and Spirtes, P.: Review of Causal Discovery Methods Based on Graphical Models, *Front. Genet.*,  
1370 10, 524, <https://doi.org/10.3389/fgene.2019.00524>, 2019.

1371 Grose, M. R., Narsey, S., Trancoso, R., Mackallah, C., Delage, F., Dowdy, A., Di Virgilio, G., Watterson, I., Dobrohotoff,  
1372 P., Rashid, H. A., Rauniyar, S., Henley, B., Thatcher, M., Syktu, J., Abramowitz, G., Evans, J. P., Su, C.-H., and Takbash,  
1373 A.: A CMIP6-based multi-model downscaling ensemble to underpin climate change services in Australia, *Clim. Serv.*, 30,  
1374 100368, <https://doi.org/10.1016/j.cliser.2023.100368>, 2023.

1375 Grundner, A., Beucler, T., Gentine, P., Iglesias-Suarez, F., Giorgetta, M. A., and Eyring, V.: Deep Learning Based Cloud  
1376 Cover Parameterization for ICON, *J. Adv. Model. Earth Syst.*, 14, <https://doi.org/10.1029/2021ms002959>, 2022.

1377 Gu, J., Wang, Z., Kuen, J., Ma, L., Shahroudy, A., Shuai, B., Liu, T., Wang, X., Wang, G., Cai, J., and Chen, T.: Recent  
1378 advances in convolutional neural networks, *Pattern Recognit.*, 77, 354–377, <https://doi.org/10.1016/j.patcog.2017.10.013>,  
1379 2018.

1380 Guendelman, I., Merlis, T. M., Cheng, K., Harris, L. M., Bretherton, C. S., Bolot, M., Zhou, L., Kaltenbaugh, A., Clark, S.  
1381 K., and Fueglistaler, S.: The Precipitation Response to Warming and CO<sub>2</sub> Increase: A Comparison of a Global Storm  
1382 Resolving Model and CMIP6 Models, *Geophys. Res. Lett.*, 51, e2023GL107008, <https://doi.org/10.1029/2023GL107008>,  
1383 2024.

- 1384 Gutowski Jr., W. J., Giorgi, F., Timbal, B., Frigon, A., Jacob, D., Kang, H.-S., Raghavan, K., Lee, B., Lennard, C., Nikulin,  
 1385 G., O'Rourke, E., Rixen, M., Solman, S., Stephenson, T., and Tangang, F.: WCRP COordinated Regional Downscaling  
 1386 EXperiment (CORDEX): a diagnostic MIP for CMIP6, *Geosci. Model Dev.*, 9, 4087–4095, [https://doi.org/10.5194/gmd-9-](https://doi.org/10.5194/gmd-9-4087-2016)  
 1387 4087-2016, 2016.
- 1388 Hagedorn, R., Doblas-Reyes, F. J., and Palmer, T. N.: The rationale behind the success of multi-model ensembles in seasonal  
 1389 forecasting – I. Basic concept, *Tellus Dyn. Meteorol. Oceanogr.*, 57, 219, <https://doi.org/10.3402/tellusa.v57i3.14657>, 2005.
- 1390 Hall, A.: Projecting regional change, *Science*, 346, 1461–1462, <https://doi.org/10.1126/science.aaa0629>, 2014.
- 1391 Hall, A., Cox, P., Huntingford, C., and Klein, S.: Progressing emergent constraints on future climate change, *Nat. Clim.*  
 1392 *Change*, 9, 269–278, <https://doi.org/10.1038/s41558-019-0436-6>, 2019.
- 1393 Hasselmann, K.: Stochastic climate models Part I. Theory, *Tellus*, 28, 473–485, <https://doi.org/10.3402/tellusa.v28i6.11316>,  
 1394 1976.
- 1395 Hawkins, E. and Sutton, R.: The Potential to Narrow Uncertainty in Regional Climate Predictions, *Bull. Am. Meteorol. Soc.*,  
 1396 90, 1095–1108, <https://doi.org/10.1175/2009BAMS2607.1>, 2009.
- 1397 Henderson, S. A., Maloney, E. D., and Son, S.-W.: Madden–Julian Oscillation Pacific Teleconnections: The Impact of the  
 1398 Basic State and MJO Representation in General Circulation Models, *J. Clim.*, 30, 4567–4587, [https://doi.org/10.1175/JCLI-](https://doi.org/10.1175/JCLI-D-16-0789.1)  
 1399 D-16-0789.1, 2017.
- 1400 Herger, N., Abramowitz, G., Knutti, R., Angélil, O., Lehmann, K., and Sanderson, B. M.: Selecting a climate model subset  
 1401 to optimise key ensemble properties, *Earth Syst. Dyn.*, 9, 135–151, <https://doi.org/10.5194/esd-9-135-2018>, 2018.
- 1402 Hilburn, K. A., Ebert-Uphoff, I., and Miller, S. D.: Development and Interpretation of a Neural-Network-Based Synthetic  
 1403 Radar Reflectivity Estimator Using GOES-R Satellite Observations, <https://doi.org/10.1175/JAMC-D-20-0084.1>, 2020.
- 1404 Hohenegger, C., Korn, P., Linardakis, L., Redler, R., Schnur, R., Adamidis, P., Bao, J., Bastin, S., Behraves, M.,  
 1405 Bergemann, M., Biercamp, J., Bockelmann, H., Brokopf, R., Brüggemann, N., Casaroli, L., Chegini, F., Datsieris, G., Esch,  
 1406 M., George, G., Giorgetta, M., Gutjahr, O., Haak, H., Hanke, M., Ilyina, T., Jahns, T., Jungclaus, J., Kern, M., Klocke, D.,  
 1407 Kluff, L., Kölling, T., Kornbluch, L., Kosukhin, S., Kroll, C., Lee, J., Mauritsen, T., Mehlmann, C., Mieslinger, T.,  
 1408 Naumann, A. K., Paccini, L., Peinado, A., Praturi, D. S., Putrasahan, D., Rast, S., Riddick, T., Roeber, N., Schmidt, H.,  
 1409 Schulzweida, U., Schütte, F., Segura, H., Shevchenko, R., Singh, V., Specht, M., Stephan, C. C., Von Storch, J.-S., Vogel,  
 1410 R., Wengel, C., Winkler, M., Ziemer, F., Marotzke, J., and Stevens, B.: ICON-Sapphire: simulating the components of the  
 1411 Earth system and their interactions at kilometer and subkilometer scales, *Geosci. Model Dev.*, 16, 779–811,  
 1412 <https://doi.org/10.5194/gmd-16-779-2023>, 2023.
- 1413 Hong, T., Wu, J., Kang, X., Yuan, M., and Duan, L.: Impacts of Different Land Use Scenarios on Future Global and  
 1414 Regional Climate Extremes, *Atmosphere*, 13, 995, <https://doi.org/10.3390/atmos13060995>, 2022.
- 1415 Iglesias-Suarez, F., Gentine, P., Solino-Fernandez, B., Beucler, T., Pritchard, M., Runge, J., and Eyring, V.: Causally-  
 1416 Informed Deep Learning to Improve Climate Models and Projections, *J. Geophys. Res. Atmospheres*, 129, e2023JD039202,  
 1417 <https://doi.org/10.1029/2023JD039202>, 2024.
- 1418 Iles, C. E., Vautard, R., Strachan, J., Joussaume, S., Eggen, B. R., and Hewitt, C. D.: The benefits of increasing resolution in  
 1419 global and regional climate simulations for European climate extremes, *Geosci. Model Dev.*, 13, 5583–5607,  
 1420 <https://doi.org/10.5194/gmd-13-5583-2020>, 2020.
- 1421 Intergovernmental Panel On Climate Change: Climate Change 2001– The Scientific Basis: Contribution of Working Group I  
 1422 to the Third Assessment Report of the Intergovernmental Panel on Climate Change, Cambridge University Press,  
 1423 Cambridge, 2001.
- 1424 Intergovernmental Panel On Climate Change: Climate Change 2007 – The Physical Science Basis: Contribution of Working  
 1425 Group I to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change., Cambridge University Press,  
 1426 Cambridge, 2007.

- 1427 Intergovernmental Panel On Climate Change (Ed.): Climate Change 2013 – The Physical Science Basis: Working Group I  
1428 Contribution to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change, 1st ed., Cambridge  
1429 University Press, <https://doi.org/10.1017/CBO9781107415324>, 2014.
- 1430 Intergovernmental Panel on Climate Change (IPCC): Climate Change 2021 – The Physical Science Basis: Working Group I  
1431 Contribution to the Sixth Assessment Report of the Intergovernmental Panel on Climate Change, Cambridge University  
1432 Press, Cambridge, <https://doi.org/10.1017/9781009157896>, 2021.
- 1433 Ivanova, D. P., Gleckler, P. J., Taylor, K. E., Durack, P. J., and Marvel, K. D.: Moving beyond the Total Sea Ice Extent in  
1434 Gauging Model Biases, *J. Clim.*, 29, 8965–8987, <https://doi.org/10.1175/JCLI-D-16-0026.1>, 2016.
- 1435 Jose, D. M., Vincent, A. M., and Dwarakish, G. S.: Improving multiple model ensemble predictions of daily precipitation  
1436 and temperature through machine learning techniques, *Sci. Rep.*, 12, 4678, <https://doi.org/10.1038/s41598-022-08786-w>,  
1437 2022.
- 1438 Joussaume, S. and Budich, R.: The Infrastructure Project of the European Network for Earth System Modelling: IS-ENES,  
1439 in: *Earth System Modelling - Volume 1*, Springer Berlin Heidelberg, Berlin, Heidelberg, 5–9, [https://doi.org/10.1007/978-3-642-36597-3\\_2](https://doi.org/10.1007/978-3-642-36597-3_2), 2013.
- 1441 Jun, M., Knutti, R., and Nychka, D. W.: Spatial Analysis to Quantify Numerical Model Bias and Dependence: How Many  
1442 Climate Models Are There?, *J. Am. Stat. Assoc.*, 103, 934–947, <https://doi.org/10.1198/016214507000001265>, 2008.
- 1443 Jung, J., Han, H., Kim, K., and Kim, H. S.: Machine Learning-Based Small Hydropower Potential Prediction under Climate  
1444 Change, *Energies*, 14, <https://doi.org/10.3390/en14123643>, 2021.
- 1445 Kaltenborn, J., Lange, C. E. E., Ramesh, V., Brouillard, P., Gurwicz, Y., Nagda, C., Runge, J., Nowack, P., and Rolnick, D.:  
1446 ClimateSet: A Large-Scale Climate Model Dataset for Machine Learning, <https://doi.org/10.48550/ARXIV.2311.03721>,  
1447 2023.
- 1448 Karmouche, S., Galytska, E., Runge, J., Meehl, G. A., Phillips, A. S., Weigel, K., and Eyring, V.: Regime-oriented causal  
1449 model evaluation of Atlantic–Pacific teleconnections in CMIP6, *Earth Syst. Dyn.*, 14, 309–344, <https://doi.org/10.5194/esd-14-309-2023>, 2023.
- 1451 Karpechko, A. Yu., Maraun, D., and Eyring, V.: Improving Antarctic Total Ozone Projections by a Process-Oriented  
1452 Multiple Diagnostic Ensemble Regression, *J. Atmospheric Sci.*, 70, 3959–3976, <https://doi.org/10.1175/JAS-D-13-071.1>,  
1453 2013.
- 1454 Katzenberger, A., Petri, S., Feulner, G., and Levermann, A.: Monsoon Planet: Bimodal Rainfall Distribution due to Barrier  
1455 Structure in Pressure Fields, *J. Clim.*, 37, 1295–1315, <https://doi.org/10.1175/JCLI-D-23-0055.1>, 2024.
- 1456 Kaufman, Z., Feldl, N., and Beaulieu, C.: Warm Arctic–Cold Eurasia pattern driven by atmospheric blocking in models and  
1457 observations, *Environ. Res. Clim.*, 3, 015006, <https://doi.org/10.1088/2752-5295/ad1f40>, 2024.
- 1458 Keenan, T. F., Luo, X., Stocker, B. D., De Kauwe, M. G., Medlyn, B. E., Prentice, I. C., Smith, N. G., Terrer, C., Wang, H.,  
1459 Zhang, Y., and Zhou, S.: A constraint on historic growth in global photosynthesis due to rising CO<sub>2</sub>, *Nat. Clim. Change*, 13,  
1460 1376–1381, <https://doi.org/10.1038/s41558-023-01867-2>, 2023.
- 1461 Kim, D., Moon, Y., Camargo, S. J., Wing, A. A., Sobel, A. H., Murakami, H., Vecchi, G. A., Zhao, M., and Page, E.:  
1462 Process-Oriented Diagnosis of Tropical Cyclones in High-Resolution GCMs, *J. Clim.*, 31, 1685–1702,  
1463 <https://doi.org/10.1175/JCLI-D-17-0269.1>, 2018.
- 1464 Kim, Y.-H., Min, S.-K., Zhang, X., Sillmann, J., and Sandstad, M.: Evaluation of the CMIP6 multi-model ensemble for  
1465 climate extreme indices, *Weather Clim. Extrem.*, 29, 100269, <https://doi.org/10.1016/j.wace.2020.100269>, 2020.
- 1466 Kirtman, B. P., Min, D., Infanti, J. M., Kinter, J. L., Paolino, D. A., Zhang, Q., Van Den Dool, H., Saha, S., Mendez, M. P.,  
1467 Becker, E., Peng, P., Tripp, P., Huang, J., DeWitt, D. G., Tippett, M. K., Barnston, A. G., Li, S., Rosati, A., Schubert, S. D.,  
1468 Rienecker, M., Suarez, M., Li, Z. E., Marshak, J., Lim, Y.-K., Tribbia, J., Pegion, K., Merryfield, W. J., Denis, B., and  
1469 Wood, E. F.: The North American Multimodel Ensemble: Phase-1 Seasonal-to-Interannual Prediction; Phase-2 toward

- 1470 Developing Intraseasonal Prediction, *Bull. Am. Meteorol. Soc.*, 95, 585–601, <https://doi.org/10.1175/BAMS-D-12-00050.1>,  
1471 2014.
- 1472 Knutson, T. R., Sirutis, J. J., Vecchi, G. A., Garner, S., Zhao, M., Kim, H.-S., Bender, M., Tuleya, R. E., Held, I. M., and  
1473 Villarini, G.: Dynamical Downscaling Projections of Twenty-First-Century Atlantic Hurricane Activity: CMIP3 and CMIP5  
1474 Model-Based Scenarios, *J. Clim.*, 26, 6591–6617, <https://doi.org/10.1175/JCLI-D-12-00539.1>, 2013.
- 1475 Knutti, R.: Should We Believe Model Predictions of Future Climate Change?, *Philos. Trans. Math. Phys. Eng. Sci.*, 366,  
1476 4647–4664, 2008.
- 1477 Knutti, R.: The end of model democracy?: An editorial comment, *Clim. Change*, 102, 395–404,  
1478 <https://doi.org/10.1007/s10584-010-9800-2>, 2010.
- 1479 Knutti, R., Furrer, R., Tebaldi, C., Cermak, J., and Meehl, G. A.: Challenges in Combining Projections from Multiple  
1480 Climate Models, <https://doi.org/10.1175/2009JCLI3361.1>, 2010a.
- 1481 Knutti, R., Abramowitz, G., Collins, M., Eyring, V., Gleckler, P. J., and Hewitson, B.: Good Practice Guidance Paper on  
1482 Assessing and Combining Multi Model Climate Projections, in: Meeting Report of the Intergovernmental Panel on Climate  
1483 Change Expert Meeting on Assessing and Combining Multi Model Climate Projections [Stocker, T.F., D. Qin, G.-K.  
1484 Plattner, M. Tignor, and P.M. Midgley (eds.)], 2010b.
- 1485 Knutti, R., Sedláček, J., Sanderson, B. M., Lorenz, R., Fischer, E. M., and Eyring, V.: A climate model projection weighting  
1486 scheme accounting for performance and interdependence, *Geophys. Res. Lett.*, 44, 1909–1918,  
1487 <https://doi.org/10.1002/2016GL072012>, 2017.
- 1488 Knutti, R., Baumberger, C., and Hirsch Hadorn, G.: Uncertainty Quantification Using Multiple Models—Prospects and  
1489 Challenges, in: *Computer Simulation Validation: Fundamental Concepts, Methodological Frameworks, and Philosophical  
1490 Perspectives*, edited by: Beisbart, C. and Saam, N. J., Springer International Publishing, Cham, 835–855,  
1491 [https://doi.org/10.1007/978-3-319-70766-2\\_34](https://doi.org/10.1007/978-3-319-70766-2_34), 2019.
- 1492 Kretschmer, M., Zappa, G., and Shepherd, T. G.: The role of Barents–Kara sea ice loss in projected polar vortex changes,  
1493 *Weather Clim. Dyn.*, 1, 715–730, <https://doi.org/10.5194/wcd-1-715-2020>, 2020.
- 1494 Krishnamurti, T. N., Kishtawal, C. M., LaRow, T. E., Bachiocchi, D. R., Zhang, Z., Williford, C. E., Gadgil, S., and  
1495 Surendran, S.: Improved Weather and Seasonal Climate Forecasts from Multimodel Superensemble, *Science*, 285, 1548–  
1496 1550, <https://doi.org/10.1126/science.285.5433.1548>, 1999.
- 1497 Kuma, P., Bender, F. A.-M., and Jönsson, A. R.: Climate Model Code Genealogy and Its Relation to Climate Feedbacks and  
1498 Sensitivity, *J. Adv. Model. Earth Syst.*, 15, e2022MS003588, <https://doi.org/10.1029/2022MS003588>, 2023.
- 1499 Kunimitsu, T., Baldissera Pacchetti, M., Ciullo, A., Sillmann, J., Shepherd, T. G., Taner, M. Ü., and van den Hurk, B.:  
1500 Representing storylines with causal networks to support decision making: Framework and example, *Clim. Risk Manag.*, 40,  
1501 100496, <https://doi.org/10.1016/j.crm.2023.100496>, 2023.
- 1502 Kyono, T., Zhang, Y., and van der Schaar, M.: CASTLE: Regularization via Auxiliary Causal Graph Discovery, in:  
1503 *Advances in Neural Information Processing Systems*, 1501–1512, 2020.
- 1504 Labe, Z. M. and Barnes, E. A.: Comparison of Climate Model Large Ensembles With Observations in the Arctic Using  
1505 Simple Neural Networks, *Earth Space Sci.*, 9, e2022EA002348, <https://doi.org/10.1029/2022EA002348>, 2022.
- 1506 Labe, Z. M., Johnson, N. C., and Delworth, T. L.: Changes in United States Summer Temperatures Revealed by Explainable  
1507 Neural Networks, *Earths Future*, 12, e2023EF003981, <https://doi.org/10.1029/2023EF003981>, 2024.
- 1508 Lambert, S. J. and Boer, G. J.: CMIP1 evaluation and intercomparison of coupled climate models, *Clim. Dyn.*, 17, 83–106,  
1509 <https://doi.org/10.1007/PL00013736>, 2001.
- 1510 LeCun, Y., Bengio, Y., and Hinton, G.: Deep learning, *Nature*, 521, 436–444, <https://doi.org/10.1038/nature14539>, 2015.
- 1511 Lehner, F. and Deser, C.: Origin, importance, and predictive limits of internal climate variability, *Environ. Res. Clim.*, 2,  
1512 023001, 2023.

- 1513 Lehner, F., Coats, S., Stocker, T. F., Pendergrass, A. G., Sanderson, B. M., Raible, C. C., and Smerdon, J. E.: Projected  
1514 drought risk in 1.5°C and 2°C warmer climates, *Geophys. Res. Lett.*, 44, 7419–7428,  
1515 <https://doi.org/10.1002/2017GL074117>, 2017.
- 1516 Lehner, F., Deser, C., Simpson, I. R., and Terray, L.: Attributing the U.S. Southwest’s recent shift into drier conditions,  
1517 *Geophys Res Lett*, 45, 6251–61, <https://doi.org/10.1029/2018GL078312>, 2018.
- 1518 Lehner, F., Deser, C., Maher, N., Marotzke, J., Fischer, E., Brunner, L., Knutti, R., and Hawkins, E.: Partitioning climate  
1519 projection uncertainty with multiple large ensembles and CMIP5/6, *Earth Syst Dyn*, 11, 491–508,  
1520 <https://doi.org/10.5194/esd-11-491-2020>, 2020.
- 1521 Li, T., Jiang, Z., Le Treut, H., Li, L., Zhao, L., and Ge, L.: Machine learning to optimize climate projection over China with  
1522 multi-model ensemble simulations, *Environ. Res. Lett.*, 16, 094028, 2021.
- 1523 Li, Y., Wu, J., Luo, J.-J., and Yang, Y. M.: Evaluating the Eastward Propagation of the MJO in CMIP5 and CMIP6 Models  
1524 Based on a Variety of Diagnostics, *J. Clim.*, 35, 1719–1743, <https://doi.org/10.1175/JCLI-D-21-0378.1>, 2022.
- 1525 Liang-Liang, L., Jian, L., and Ru-Cong, Y.: Evaluation of CMIP6 HighResMIP models in simulating precipitation over  
1526 Central Asia, *Adv. Clim. Change Res.*, 13, 1–13, <https://doi.org/10.1016/j.accre.2021.09.009>, 2022.
- 1527 Lin, X., Zhen, H.-L., Li, Z., Zhang, Q.-F., and Kwong, S.: Pareto Multi-Task Learning, in: *Advances in Neural Information*  
1528 *Processing Systems*, 2019.
- 1529 Liu, Y., Fan, K., Chen, L., Ren, H.-L., Wu, Y., and Liu, C.: An operational statistical downscaling prediction model of the  
1530 winter monthly temperature over China based on a multi-model ensemble, *Atmospheric Res.*, 249, 105262,  
1531 <https://doi.org/10.1016/j.atmosres.2020.105262>, 2021.
- 1532 Lovenduski, N. S., McKinley, G. A., Fay, A. R., Lindsay, K., and Long, M. C.: Partitioning uncertainty in ocean carbon  
1533 uptake projections: internal variability, emission scenario, and model structure, *Glob Biogeochem Cycles*, 30, 1276–87,  
1534 <https://doi.org/10.1002/2016GB005426>, 2016.
- 1535 Lu, D. and Ricciuto, D.: Efficient surrogate modeling methods for large-scale Earth system models based on machine-  
1536 learning techniques, *Geosci. Model Dev.*, 12, 1791–1807, <https://doi.org/10.5194/gmd-12-1791-2019>, 2019.
- 1537 Luo, Y., Peng, J., and Ma, J.: When causal inference meets deep learning, *Nat. Mach. Intell.*, 2, 426–427,  
1538 <https://doi.org/10.1038/s42256-020-0218-x>, 2020.
- 1539 Maher, N., Phillips, A. S., Deser, C., Wills, R. C. J., Lehner, F., Fasullo, J., Caron, J. M., Brunner, L., and Beyerle, U.: The  
1540 updated Multi-Model Large Ensemble Archive and the Climate Variability Diagnostics Package: New tools for the study of  
1541 climate variability and change, <https://doi.org/10.5194/egusphere-2024-3684>, 19 December 2024.
- 1542 Maher, N., Phillips, A. S., Deser, C., Wills, R. C. J., Lehner, F., Fasullo, J., Caron, J. M., Brunner, L., Beyerle, U., and  
1543 Jeffree, J.: The updated Multi-Model Large Ensemble Archive and the Climate Variability Diagnostics Package: new tools  
1544 for the study of climate variability and change, *Geosci. Model Dev.*, 18, 6341–6365, <https://doi.org/10.5194/gmd-18-6341-2025>, 2025.
- 1546 Maloney, E. D., Gettelman, A., Ming, Y., Neelin, J. D., Barrie, D., Mariotti, A., Chen, C.-C., Coleman, D. R. B., Kuo, Y.-H.,  
1547 Singh, B., Annamalai, H., Berg, A., Booth, J. F., Camargo, S. J., Dai, A., Gonzalez, A., Hafner, J., Jiang, X., Jing, X., Kim,  
1548 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.:  
1549 Process-Oriented Evaluation of Climate and Weather Forecasting Models, *Bull. Am. Meteorol. Soc.*, 100, 1665–1686,  
1550 <https://doi.org/10.1175/BAMS-D-18-0042.1>, 2019.
- 1551 Manabe, S. and Bryan, K.: Climate Calculations with a Combined Ocean-Atmosphere Model, *J. Atmospheric Sci.*, 26, 786–  
1552 789, [https://doi.org/10.1175/1520-0469\(1969\)026%253C0786:CCWACO%253E2.0.CO;2](https://doi.org/10.1175/1520-0469(1969)026%253C0786:CCWACO%253E2.0.CO;2), 1969.
- 1553 Manabe, S. and Strickler, R. F.: Thermal Equilibrium of the Atmosphere with a Convective Adjustment, *J. Atmospheric Sci.*,  
1554 21, 361–385, [https://doi.org/10.1175/1520-0469\(1964\)021%253C0361:TEOTAW%253E2.0.CO;2](https://doi.org/10.1175/1520-0469(1964)021%253C0361:TEOTAW%253E2.0.CO;2), 1964.
- 1555 Manabe, S. and Wetherald, R. T.: Thermal Equilibrium of the Atmosphere with a Given Distribution of Relative Humidity,

- 1556 J. Atmospheric Sci., 24, 241–259, [https://doi.org/10.1175/1520-0469\(1967\)024%253C0241:TEOTAW%253E2.0.CO;2](https://doi.org/10.1175/1520-0469(1967)024%253C0241:TEOTAW%253E2.0.CO;2),  
1557 1967.
- 1558 Mankin, J. S. and Diffenbaugh, N. S.: Influence of temperature and precipitation variability on near-term snow trends, *Clim.*  
1559 *Dyn.*, 45, 1099–1116, <https://doi.org/10.1007/s00382-014-2357-4>, 2015.
- 1560 Maraun, D.: Bias Correcting Climate Change Simulations - a Critical Review, *Curr. Clim. Change Rep.*, 2, 211–220,  
1561 <https://doi.org/10.1007/s40641-016-0050-x>, 2016.
- 1562 Maraun, D., Shepherd, T. G., Widmann, M., Zappa, G., Walton, D., Gutiérrez, J. M., Hagemann, S., Richter, I., Soares, P.  
1563 M. M., Hall, A., and Mearns, L. O.: Towards process-informed bias correction of climate change simulations, *Nat. Clim.*  
1564 *Change*, 7, 764–773, <https://doi.org/10.1038/nclimate3418>, 2017.
- 1565 Marotzke, J.: Quantifying the irreducible uncertainty in near-term climate projections, *WIREs Clim. Change*, 10, e563,  
1566 <https://doi.org/10.1002/wcc.563>, 2019.
- 1567 Masson, D. and Knutti, R.: Climate model genealogy, *Geophys. Res. Lett.*, 38, <https://doi.org/10.1029/2011GL046864>,  
1568 2011.
- 1569 Mauritsen, T., Bader, J., Becker, T., Behrens, J., Bittner, M., Brokopf, R., Brovkin, V., Claussen, M., Crueger, T., Esch, M.,  
1570 Fast, I., Fiedler, S., Fläschner, D., Gayler, V., Giorgetta, M., Goll, D. S., Haak, H., Hagemann, S., Hedemann, C.,  
1571 Hohenegger, C., Ilyina, T., Jahns, T., Jimenéz-de-la-Cuesta, D., Jungclaus, J., Kleinen, T., Kloster, S., Kracher, D., Kinne,  
1572 S., Kleberg, D., Lasslop, G., Kornblueh, L., Marotzke, J., Matei, D., Meraner, K., Mikolajewicz, U., Modali, K., Möbis, B.,  
1573 Müller, W. A., Nabel, J. E. M. S., Nam, C. C. W., Notz, D., Nyawira, S., Paulsen, H., Peters, K., Pincus, R., Pohlmann, H.,  
1574 Pongratz, J., Popp, M., Raddatz, T. J., Rast, S., Redler, R., Reick, C. H., Rohrschneider, T., Schemann, V., Schmidt, H.,  
1575 Schnur, R., Schulzweida, U., Six, K. D., Stein, L., Stemmler, I., Stevens, B., Von Storch, J., Tian, F., Voigt, A., Vrese, P.,  
1576 Wieners, K., Wilkenskield, S., Winkler, A., and Roeckner, E.: Developments in the MPI-M Earth System Model version 1.2  
1577 (MPI-ESM1.2) and Its Response to Increasing CO<sub>2</sub>, *J. Adv. Model. Earth Syst.*, 11, 998–1038,  
1578 <https://doi.org/10.1029/2018MS001400>, 2019.
- 1579 Meehl, G. A., Boer, G. J., Covey, C., Latif, M., and Stouffer, R. J.: The Coupled Model Intercomparison Project (CMIP),  
1580 *Bull. Am. Meteorol. Soc.*, 81, 313–318, 2000.
- 1581 Mendlik, T. and Gobiet, A.: Selecting climate simulations for impact studies based on multivariate patterns of climate  
1582 change, *Clim. Change*, 135, 381–393, <https://doi.org/10.1007/s10584-015-1582-0>, 2016.
- 1583 Merlis, T. M., Cheng, K.-Y., Guendelman, I., Harris, L., Bretherton, C. S., Bolot, M., Zhou, L., Kaltenbaugh, A., Clark, S.  
1584 K., Vecchi, G. A., and Fueglistaler, S.: Climate sensitivity and relative humidity changes in global storm-resolving model  
1585 simulations of climate change, *Sci. Adv.*, 10, eadn5217, <https://doi.org/10.1126/sciadv.adn5217>, 2024.
- 1586 Merrifield, A. L., Brunner, L., Lorenz, R., Medhaug, I., and Knutti, R.: An investigation of weighting schemes suitable for  
1587 incorporating large ensembles into multi-model ensembles, *Earth Syst. Dyn.*, 11, 807–834, <https://doi.org/10.5194/esd-11-1588>  
1588 807-2020, 2020.
- 1589 Merrifield, A. L., Brunner, L., Lorenz, R., Humphrey, V., and Knutti, R.: Climate model Selection by Independence,  
1590 Performance, and Spread (ClimSIPS v1.0.1) for regional applications, *Geosci. Model Dev.*, 16, 4715–4747,  
1591 <https://doi.org/10.5194/gmd-16-4715-2023>, 2023.
- 1592 Milinski, S., Maher, N., and Olonscheck, D.: How large does a large ensemble need to be?, *Earth Syst. Dyn.*, 11, 885–901,  
1593 <https://doi.org/10.5194/esd-11-885-2020>, 2020.
- 1594 Min, Y., Lim, C., Yoo, J., Kim, H., Kryjov, V. N., Jeong, D., Lim, A., Ham, S., Chen, M., Xiao, Y., Gagnon, N., Muncaster,  
1595 R., Liu, P., Borrelli, A., Ji, H., Lee, J., Jo, S., Kiktev, D., Tolstykh, M., Matyugin, V., McLean, P., and Molod, A. M.: A  
1596 Diachronic Assessment of Advances in Seasonal Forecasting: Evolution of the APCC Multi-Model Ensemble Prediction  
1597 System Over the Last Two Decades, *Geophys. Res. Lett.*, 52, e2025GL116416, <https://doi.org/10.1029/2025GL116416>,  
1598 2025.
- 1599 Moon, Y., Kim, D., Camargo, S. J., Wing, A. A., Sobel, A. H., Murakami, H., Reed, K. A., Scoccimarro, E., Vecchi, G. A.,

1600 Wehner, M. F., Zarzycki, C. M., and Zhao, M.: Azimuthally Averaged Wind and Thermodynamic Structures of Tropical  
1601 Cyclones in Global Climate Models and Their Sensitivity to Horizontal Resolution, *J. Clim.*, 33, 1575–1595,  
1602 <https://doi.org/10.1175/JCLI-D-19-0172.1>, 2020.

1603 Mudryk, L., Santolaria-Otín, M., Krinner, G., Ménégos, M., Derksen, C., Brutel-Vuilmet, C., Brady, M., and Essery, R.:  
1604 Historical Northern Hemisphere snow cover trends and projected changes in the CMIP6 multi-model ensemble, *The*  
1605 *Cryosphere*, 14, 2495–2514, <https://doi.org/10.5194/tc-14-2495-2020>, 2020.

1606 Nam, C., Bony, S., Dufresne, J. -L., and Chepfer, H.: The ‘too few, too bright’ tropical low-cloud problem in CMIP5  
1607 models, *Geophys. Res. Lett.*, 39, 2012GL053421, <https://doi.org/10.1029/2012GL053421>, 2012.

1608 The Climate Data Guide: Regridding Overview: <https://climatedataguide.ucar.edu/climate-tools/regridding-overview>.

1609 Neelin, J. D., Krasting, J. P., Radhakrishnan, A., Liptak, J., Jackson, T., Ming, Y., Dong, W., Gettelman, A., Coleman, D. R.,  
1610 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.:  
1611 Process-Oriented Diagnostics: Principles, Practice, Community Development, and Common Standards, *Bull. Am. Meteorol.*  
1612 *Soc.*, 104, E1452–E1468, <https://doi.org/10.1175/BAMS-D-21-0268.1>, 2023.

1613 Nijse, F. J. M. M., Cox, P. M., and Williamson, M. S.: Emergent constraints on transient climate response (TCR) and  
1614 equilibrium climate sensitivity (ECS) from historical warming in CMIP5 and CMIP6 models, *Earth Syst. Dyn.*, 11, 737–750,  
1615 <https://doi.org/10.5194/esd-11-737-2020>, 2020.

1616 Nolan, P. and Flanagan, J.: High-resolution climate projections for Ireland - a multi-model ensemble approach: 2014-CCRP-  
1617 MS.23, Online version., Environmental Protection Agency, Johnstown Castle, Co. Wexford, Ireland, 1 pp., 2020.

1618 Notz, D., Jahn, A., Holland, M., Hunke, E., Massonnet, F., Stroeve, J., Tremblay, B., and Vancoppenolle, M.: The CMIP6  
1619 Sea-Ice Model Intercomparison Project (SIMIP): understanding sea ice through climate-model simulations, *Geosci. Model*  
1620 *Dev.*, 9, 3427–3446, <https://doi.org/10.5194/gmd-9-3427-2016>, 2016.

1621 Nowack, P., Runge, J., Eyring, V., and Haigh, J. D.: Causal networks for climate model evaluation and constrained  
1622 projections, *Nat. Commun.*, 11, 1415, <https://doi.org/10.1038/s41467-020-15195-y>, 2020.

1623 Nwokolo, S. C., Obiwulu, A. U., and Ogbulezie, J. C.: Machine learning and analytical model hybridization to assess the  
1624 impact of climate change on solar PV energy production, *Phys. Chem. Earth Parts ABC*, 130, 103389,  
1625 <https://doi.org/10.1016/j.pce.2023.103389>, 2023.

1626 Olonscheck, D. and Notz, D.: Consistently Estimating Internal Climate Variability from Climate Model Simulations, *J.*  
1627 *Clim.*, 30, 9555–9573, <https://doi.org/10.1175/JCLI-D-16-0428.1>, 2017.

1628 O’Neill, B. C., Kriegler, E., Riahi, K., Ebi, K. L., Hallegatte, S., Carter, T. R., Mathur, R., and van Vuuren, D. P.: A new  
1629 scenario framework for climate change research: the concept of shared socioeconomic pathways, *Clim. Change*, 122, 387–  
1630 400, <https://doi.org/10.1007/s10584-013-0905-2>, 2014.

1631 O’Neill, B. C., Kriegler, E., Ebi, K. L., Kemp-Benedict, E., Riahi, K., Rothman, D. S., Van Ruijven, B. J., Van Vuuren, D.  
1632 P., Birkmann, J., Kok, K., Levy, M., and Solecki, W.: The roads ahead: Narratives for shared socioeconomic pathways  
1633 describing world futures in the 21st century, *Glob. Environ. Change*, 42, 169–180,  
1634 <https://doi.org/10.1016/j.gloenvcha.2015.01.004>, 2017.

1635 Oxarart, A. and Parker, L.: Global Climate Models and Land Management, USDA California Climate Hub, 2024.

1636 Palmer, T. E., McSweeney, C. F., Booth, B. B. B., Priestley, M. D. K., Davini, P., Brunner, L., Borchert, L., and Menary, M.  
1637 B.: Performance-based sub-selection of CMIP6 models for impact assessments in Europe, *Earth Syst. Dyn.*, 14, 457–483,  
1638 <https://doi.org/10.5194/esd-14-457-2023>, 2023.

1639 Palmer, T. n, Doblas-Reyes, F. j, Hagedorn, R., and Weisheimer, A.: Probabilistic prediction of climate using multi-model  
1640 ensembles: from basics to applications, *Philos. Trans. R. Soc. B Biol. Sci.*, 360, 1991–1998,  
1641 <https://doi.org/10.1098/rstb.2005.1750>, 2005.

1642 Pennell, C. and Reichler, T.: On the Effective Number of Climate Models, <https://doi.org/10.1175/2010JCLI3814.1>, 2011.

- 1643 Pérez-Carrasquilla, J. S., Molina, M. J., Mayer, K. J., Dagon, K., Fasullo, J. T., & Simpson, I. R.: Observed and modeled  
1644 amplification of the frequency, duration, and extreme heat impacts of the Pacific trough regime. *Earth's Future*, 13(12),  
1645 e2025EF007140, 2025.
- 1646 Phillips, A., Deser, C., Fasullo, J., Schneider, D. P., and Simpson, I. R.: Assessing Climate Variability and Change in Model  
1647 Large Ensembles: A User's Guide to the "Climate Variability Diagnostics Package for Large Ensembles,"  
1648 <https://doi.org/10.5065/H7C7-F961>, 2020.
- 1649 Phillips, A. S., Deser, C., and Fasullo, J.: Evaluating Modes of Variability in Climate Models, *Eos Trans. Am. Geophys.*  
1650 *Union*, 95, 453–455, <https://doi.org/10.1002/2014EO490002>, 2014.
- 1651 Phillips, T. J. and Gleckler, P. J.: Evaluation of continental precipitation in 20th century climate simulations: The utility of  
1652 multimodel statistics, *Water Resour. Res.*, 42, 2005WR004313, <https://doi.org/10.1029/2005WR004313>, 2006.
- 1653 Pichelli, E., Coppola, E., Sobolowski, S., Ban, N., Giorgi, F., Stocchi, P., Alias, A., Belušić, D., Berthou, S., Caillaud, C.,  
1654 Cardoso, R. M., Chan, S., Christensen, O. B., Dobler, A., de Vries, H., Goergen, K., Kendon, E. J., Keuler, K., Lenderink,  
1655 G., Lorenz, T., Mishra, A. N., Panitz, H.-J., Schär, C., Soares, P. M. M., Truhetz, H., and Vergara-Temprado, J.: The first  
1656 multi-model ensemble of regional climate simulations at kilometer-scale resolution part 2: historical and future simulations  
1657 of precipitation, *Clim. Dyn.*, 56, 3581–3602, <https://doi.org/10.1007/s00382-021-05657-4>, 2021.
- 1658 Pincus, R., Barker, H. W., and Morcrette, J.: A fast, flexible, approximate technique for computing radiative transfer in  
1659 inhomogeneous cloud fields, *J. Geophys. Res. Atmospheres*, 108, 2002JD003322, <https://doi.org/10.1029/2002JD003322>,  
1660 2003.
- 1661 Pincus, R., Batstone, C. P., Hofmann, R. J. P., Taylor, K. E., and Glecker, P. J.: Evaluating the present-day simulation of  
1662 clouds, precipitation, and radiation in climate models, *J. Geophys. Res. Atmospheres*, 113,  
1663 <https://doi.org/10.1029/2007JD009334>, 2008.
- 1664 Planton, Y. Y., Guilyardi, E., Wittenberg, A. T., Lee, J., Gleckler, P. J., Bayr, T., McGregor, S., McPhaden, M. J., Power, S.,  
1665 Roehrig, R., Vialard, J., and Voltaire, A.: Evaluating Climate Models with the CLIVAR 2020 ENSO Metrics Package, *Bull.*  
1666 *Am. Meteorol. Soc.*, 102, E193–E217, <https://doi.org/10.1175/BAMS-D-19-0337.1>, 2021.
- 1667 Polkova\*, I., Afargan-Gerstman, H., Domeisen, D. I. V., King, M. P., Ruggieri, P., Athanasiadis, P., Dobrynin, M., Aarnes,  
1668 Ø., Kretschmer, M., and Baehr, J.: Predictors and prediction skill for marine cold-air outbreaks over the Barents Sea, *Q. J. R.*  
1669 *Meteorol. Soc.*, 147, 2638–2656, <https://doi.org/10.1002/qj.4038>, 2021.
- 1670 Quesada, B., Arneeth, A., and de Noblet-Ducoudré, N.: Atmospheric, radiative, and hydrologic effects of future land use and  
1671 land cover changes: A global and multimodel climate picture, *J. Geophys. Res. Atmospheres*, 122, 5113–5131,  
1672 <https://doi.org/10.1002/2016JD025448>, 2017.
- 1673 Rackow, T., Pedruzo-Bagazgoitia, X., Becker, T., Milinski, S., Sandu, I., Aguridan, R., Bechtold, P., Beyer, S., Bidlot, J.,  
1674 Boussetta, S., Deconinck, W., Diamantakis, M., Dueben, P., Dutra, E., Forbes, R., Ghosh, R., Goessling, H. F., Hadade, I.,  
1675 Hegewald, J., Jung, T., Keeley, S., Kluft, L., Koldunov, N., Koldunov, A., Kölling, T., Kousal, J., Kühnlein, C., Maciel, P.,  
1676 Mogensen, K., Quintino, T., Polichtchouk, I., Reuter, B., Sármany, D., Scholz, P., Sidorenko, D., Streffing, J., Sützl, B.,  
1677 Takasuka, D., Tietsche, S., Valentini, M., Vannière, B., Wedi, N., Zampieri, L., and Ziemann, F.: Multi-year simulations at  
1678 kilometre scale with the Integrated Forecasting System coupled to FESOM2.5 and NEMOV3.4, *Geosci. Model Dev.*, 18, 33–  
1679 69, <https://doi.org/10.5194/gmd-18-33-2025>, 2025.
- 1680 Rader, J. K., Barnes, E. A., Ebert-Uphoff, I., and Anderson, C.: Detection of Forced Change Within Combined Climate  
1681 Fields Using Explainable Neural Networks, *J. Adv. Model. Earth Syst.*, 14, e2021MS002941,  
1682 <https://doi.org/10.1029/2021MS002941>, 2022.
- 1683 Räisänen, J.: Objective comparison of patterns of CO<sub>2</sub> induced climate change in coupled GCM experiments, *Clim. Dyn.*,  
1684 13, 197–211, <https://doi.org/10.1007/s003820050160>, 1997.
- 1685 Räisänen, J. and Palmer, T. N.: A Probability and Decision-Model Analysis of a Multimodel Ensemble of Climate Change  
1686 Simulations, *J. Clim.*, 14, 3212–3226, [https://doi.org/10.1175/1520-0442\(2001\)014%253C3212:APADMA%253E2.0.CO;2](https://doi.org/10.1175/1520-0442(2001)014%253C3212:APADMA%253E2.0.CO;2),

- 1687 2001.
- 1688 Rampal, N., Hobeichi, S., Gibson, P. B., Baño-Medina, J., Abramowitz, G., Beucler, T., González-Abad, J., Chapman, W.,  
1689 Harder, P., and Gutiérrez, J. M.: Enhancing Regional Climate Downscaling through Advances in Machine Learning, *Artif.*  
1690 *Intell. Earth Syst.*, 3, 230066, <https://doi.org/10.1175/AIES-D-23-0066.1>, 2024.
- 1691 Rasp, S., Pritchard, M. S., and Gentine, P.: Deep learning to represent subgrid processes in climate models, *Proc. Natl. Acad.*  
1692 *Sci.*, 115, 9684–9689, <https://doi.org/10.1073/pnas.1810286115>, 2018.
- 1693 Reichler, T. and Kim, J.: How Well Do Coupled Models Simulate Today’s Climate?, [https://doi.org/10.1175/BAMS-89-3-](https://doi.org/10.1175/BAMS-89-3-303)  
1694 303, 2008.
- 1695 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., and Prabhat: Deep learning and process  
1696 understanding for data-driven Earth system science, *Nature*, 566, 195–204, <https://doi.org/10.1038/s41586-019-0912-1>,  
1697 2019.
- 1698 Riahi, K., van Vuuren, D. P., Kriegler, E., Edmonds, J., O’Neill, B. C., Fujimori, S., Bauer, N., Calvin, K., Dellink, R.,  
1699 Fricko, O., Lutz, W., Popp, A., Cuaresma, J. C., Kc, S., Leimbach, M., Jiang, L., Kram, T., Rao, S., Emmerling, J., Ebi, K.,  
1700 Hasegawa, T., Havlik, P., Humpenöder, F., Da Silva, L. A., Smith, S., Stehfest, E., Bosetti, V., Eom, J., Gernaat, D., Masui,  
1701 T., Rogelj, J., Strefler, J., Drouet, L., Krey, V., Luderer, G., Harmsen, M., Takahashi, K., Baumstark, L., Doelman, J. C.,  
1702 Kainuma, M., Klimont, Z., Marangoni, G., Lotze-Campen, H., Obersteiner, M., Tabeau, A., and Tavoni, M.: The Shared  
1703 Socioeconomic Pathways and their energy, land use, and greenhouse gas emissions implications: An overview, *Glob.*  
1704 *Environ. Change*, 42, 153–168, <https://doi.org/10.1016/j.gloenvcha.2016.05.009>, 2017.
- 1705 Ricard, L., Falasca, F., Runge, J., and Nenes, A.: network-based constraint to evaluate climate sensitivity, *Nat. Commun.*,  
1706 15, 6942, <https://doi.org/10.1038/s41467-024-50813-z>, 2024.
- 1707 Roach, L. A., Dean, S. M., and Renwick, J. A.: Consistent biases in Antarctic sea ice concentration simulated by climate  
1708 models, *The Cryosphere*, 12, 365–383, <https://doi.org/10.5194/tc-12-365-2018>, 2018.
- 1709 Rodgers, K. B., Lin, J., and Frölicher, T. L.: Emergence of multiple ocean ecosystem drivers in a large ensemble suite with  
1710 an Earth system model, *Biogeosciences*, 12, 3301–20, <https://doi.org/10.5194/bg-12-3301-2015>, 2015.
- 1711 Rojpratak, S. and Supharatid, S.: Regional extreme precipitation index: Evaluations and projections from the multi-model  
1712 ensemble CMIP5 over Thailand, *Weather Clim. Extrem.*, 37, 100475, <https://doi.org/10.1016/j.wace.2022.100475>, 2022.
- 1713 Roy, I., Gagnon, A. S., and Siingh, D.: Evaluating ENSO teleconnections using observations and CMIP5 models, *Theor.*  
1714 *Appl. Climatol.*, 136, 1085–1098, <https://doi.org/10.1007/s00704-018-2536-z>, 2018.
- 1715 Roy, I., Tedeschi, R. G., and Collins, M.: ENSO teleconnections to the Indian summer monsoon under changing climate, *Int.*  
1716 *J. Climatol.*, 39, 3031–3042, <https://doi.org/10.1002/joc.5999>, 2019.
- 1717 Runge, J., Bathiany, S., Bollt, E., Camps-Valls, G., Coumou, D., Deyle, E., Glymour, C., Kretschmer, M., Mahecha, M. D.,  
1718 Muñoz-Marí, J., Van Nes, E. H., Peters, J., Quax, R., Reichstein, M., Scheffer, M., Schölkopf, B., Spirtes, P., Sugihara, G.,  
1719 Sun, J., Zhang, K., and Zscheischler, J.: Inferring causation from time series in Earth system sciences, *Nat. Commun.*, 10,  
1720 2553, <https://doi.org/10.1038/s41467-019-10105-3>, 2019.
- 1721 Runge, J., Tibau, X.-A., Bruhns, M., Muñoz-Marí, J., and Camps-Valls, G.: The Causality for Climate Competition, in:  
1722 *Proceedings of the NeurIPS 2019 Competition and Demonstration Track*, 110–120, 2020.
- 1723 Runge, J., Gerhardus, A., Varando, G., Eyring, V., and Camps-Valls, G.: Causal inference for time series, *Nat. Rev. Earth*  
1724 *Environ.*, 4, 487–505, <https://doi.org/10.1038/s43017-023-00431-y>, 2023.
- 1725 Rupe, A., Crutchfield, J. P., Kashinath, K., and Prabhat: A Physics-Based Approach to Unsupervised Discovery of Coherent  
1726 Structures in Spatiotemporal Systems, <https://doi.org/10.48550/ARXIV.1709.03184>, 2017.
- 1727 Russo, F. and Toni, F.: Causal Discovery and Knowledge Injection for Contestable Neural Networks (with Appendices),  
1728 <https://doi.org/10.48550/ARXIV.2205.09787>, 2022.
- 1729 Rypkema, D. and Tuljapurkar, S.: Modeling extreme climatic events using the generalized extreme value (GEV) distribution,

- 1730 in: *Handbook of Statistics*, vol. 44, Elsevier, 39–71, <https://doi.org/10.1016/bs.host.2020.12.002>, 2021.
- 1731 Sachindra, D. A., Ahmed, K., Rashid, Md. M., Shahid, S., and Perera, B. J. C.: Statistical downscaling of precipitation using  
1732 machine learning techniques, *Atmospheric Res.*, 212, 240–258, <https://doi.org/10.1016/j.atmosres.2018.05.022>, 2018.
- 1733 Sanderson, B. M. and Knutti, R.: On the interpretation of constrained climate model ensembles, *Geophys. Res. Lett.*, 39,  
1734 <https://doi.org/10.1029/2012GL052665>, 2012.
- 1735 Sanderson, B. M., Knutti, R., Aina, T., Christensen, C., Faull, N., Frame, D. J., Ingram, W. J., Piani, C., Stainforth, D. A.,  
1736 Stone, D. A., and Allen, M. R.: Constraints on Model Response to Greenhouse Gas Forcing and the Role of Subgrid-Scale  
1737 Processes, *J. Clim.*, 21, 2384–2400, <https://doi.org/10.1175/2008JCLI1869.1>, 2008.
- 1738 Sanderson, B. M., Knutti, R., and Caldwell, P.: A Representative Democracy to Reduce Interdependency in a Multimodel  
1739 Ensemble, <https://doi.org/10.1175/JCLI-D-14-00362.1>, 2015.
- 1740 Sanderson, B. M., Pendergrass, A. G., Koven, C. D., Brient, F., Booth, B. B. B., Fisher, R. A., and Knutti, R.: The potential  
1741 for structural errors in emergent constraints, *Earth Syst. Dyn.*, 12, 899–918, <https://doi.org/10.5194/esd-12-899-2021>, 2021.
- 1742 Santer, B. D., Thorne, P. W., Haimberger, L., Taylor, K. E., Wigley, T. M. L., Lanzante, J. R., Solomon, S., Free, M.,  
1743 Gleckler, P. J., Jones, P. D., Karl, T. R., Klein, S. A., Mears, C., Nychka, D., Schmidt, G. A., Sherwood, S. C., and Wentz, F.  
1744 J.: Consistency of modelled and observed temperature trends in the tropical troposphere, *Int. J. Climatol.*, 28, 1703–1722,  
1745 <https://doi.org/10.1002/joc.1756>, 2008.
- 1746 Santer, B. D., Taylor, K. E., Gleckler, P. J., Bonfils, C., Barnett, T. P., Pierce, D. W., Wigley, T. M. L., Mears, C., Wentz, F.  
1747 J., Brüggemann, W., Gillett, N. P., Klein, S. A., Solomon, S., Stott, P. A., and Wehner, M. F.: Incorporating model quality  
1748 information in climate change detection and attribution studies, *Proc. Natl. Acad. Sci.*, 106, 14778–14783,  
1749 <https://doi.org/10.1073/pnas.0901736106>, 2009.
- 1750 Schär, C., Fuhrer, O., Arteaga, A., Ban, N., Charpilloz, C., Di Girolamo, S., Hentgen, L., Hoefler, T., Lapillonne, X.,  
1751 Leutwyler, D., Osterried, K., Panosetti, D., Rüdüsühli, S., Schlemmer, L., Schulthess, T. C., Sprenger, M., Ubbiali, S., and  
1752 Wernli, H.: Kilometer-Scale Climate Models: Prospects and Challenges, *Bull. Am. Meteorol. Soc.*, 101, E567–E587,  
1753 <https://doi.org/10.1175/BAMS-D-18-0167.1>, 2020.
- 1754 Scher, S.: Toward Data-Driven Weather and Climate Forecasting: Approximating a Simple General Circulation Model With  
1755 Deep Learning, *Geophys. Res. Lett.*, 45, 12,616–12,622, <https://doi.org/10.1029/2018GL080704>, 2018.
- 1756 Schlunegger, S., Rodgers, K. B., Sarmiento, J. L., Frölicher, T. L., Dunne, J. P., Ishii, M., and Slater, R.: Emergence of  
1757 anthropogenic signals in the ocean carbon cycle, *Nat. Clim. Change*, 9, 719–725, [https://doi.org/10.1038/s41558-019-0553-](https://doi.org/10.1038/s41558-019-0553-2)  
1758 2, 2019.
- 1759 Schneider, T., Teixeira, J., Bretherton, C. S., Brient, F., Pressel, K. G., Schär, C., and Siebesma, A. P.: Climate goals and  
1760 computing the future of clouds, *Nat. Clim. Change*, 7, 3–5, <https://doi.org/10.1038/nclimate3190>, 2017.
- 1761 Scholkopf, B., Locatello, F., Bauer, S., Ke, N. R., Kalchbrenner, N., Goyal, A., and Bengio, Y.: Toward Causal  
1762 Representation Learning, *Proc. IEEE*, 109, 612–634, <https://doi.org/10.1109/jproc.2021.3058954>, 2021.
- 1763 Sener, O. and Koltun, V.: Multi-Task Learning as Multi-Objective Optimization, in: *Advances in Neural Information*  
1764 *Processing Systems*, 2018.
- 1765 Seneviratne, S. I., Nicholls, N., Easterling, D., Goodess, C. M., Kanae, S., Kossin, J., Luo, Y., Marengo, J., McInnes, K.,  
1766 Rahimi, M., Reichstein, M., Sorteberg, A., Vera, C., Zhang, X., Rusticucci, M., Semenov, V., Alexander, L. V., Allen, S.,  
1767 Benito, G., Cavazos, T., Clague, J., Conway, D., Della-Marta, P. M., Gerber, M., Gong, S., Goswami, B. N., Hemer, M.,  
1768 Huggel, C., Van Den Hurk, B., Kharin, V. V., Kitoh, A., Tank, A. M. G. K., Li, G., Mason, S., McGuire, W., Van  
1769 Oldenborgh, G. J., Orłowsky, B., Smith, S., Thiaw, W., Velegarakis, A., Yiou, P., Zhang, T., Zhou, T., and Zwiers, F. W.:  
1770 Changes in Climate Extremes and their Impacts on the Natural Physical Environment, in: *Managing the Risks of Extreme*  
1771 *Events and Disasters to Advance Climate Change Adaptation*, edited by: Field, C. B., Barros, V., Stocker, T. F., and Dahe,  
1772 Q., Cambridge University Press, 109–230, <https://doi.org/10.1017/CBO9781139177245.006>, 2012.

- 1773 Sexton, D. M. H., McSweeney, C. F., Rostron, J. W., Yamazaki, K., Booth, B. B. B., Murphy, J. M., Regayre, L., Johnson, J.  
1774 S., and Karmalkar, A. V.: A perturbed parameter ensemble of HadGEM3-GC3.05 coupled model projections: part 1:  
1775 selecting the parameter combinations, *Clim. Dyn.*, 56, 3395–3436, <https://doi.org/10.1007/s00382-021-05709-9>, 2021.
- 1776 Shaw, T. A., Arblaster, J. M., Birner, T., Butler, A. H., Domeisen, D. I. V., Garfinkel, C. I., Garny, H., Grise, K. M., and  
1777 Karpechko, A. Yu.: Emerging Climate Change Signals in Atmospheric Circulation, *AGU Adv.*, 5, e2024AV001297,  
1778 <https://doi.org/10.1029/2024AV001297>, 2024.
- 1779 Shepherd, T. G.: Atmospheric circulation as a source of uncertainty in climate change projections, *Nat. Geosci.*, 7, 703–708,  
1780 <https://doi.org/10.1038/ngeo2253>, 2014.
- 1781 Shepherd, T. G.: Storyline approach to the construction of regional climate change information, *Proc. R. Soc. Math. Phys.*  
1782 *Eng. Sci.*, 475, 20190013, <https://doi.org/10.1098/rspa.2019.0013>, 2019.
- 1783 Shepherd, T. G., Boyd, E., Calel, R. A., Chapman, S. C., Dessai, S., Dima-West, I. M., Fowler, H. J., James, R., Maraun, D.,  
1784 Martius, O., Senior, C. A., Sobel, A. H., Stainforth, D. A., Tett, S. F. B., Trenberth, K. E., Van Den Hurk, B. J. J. M.,  
1785 Watkins, N. W., Wilby, R. L., and Zenghelis, D. A.: Storylines: an alternative approach to representing uncertainty in  
1786 physical aspects of climate change, *Clim. Change*, 151, 555–571, <https://doi.org/10.1007/s10584-018-2317-9>, 2018.
- 1787 Shetty, S., Umesh, P., and Shetty, A.: The effectiveness of machine learning-based multi-model ensemble predictions of  
1788 CMIP6 in Western Ghats of India, *Int. J. Climatol.*, 43, 5029–5054, <https://doi.org/10.1002/joc.8131>, 2023.
- 1789 Shin, Y., Lee, Y., and Park, J.-S.: A Weighting Scheme in A Multi-Model Ensemble for Bias-Corrected Climate Simulation,  
1790 *Atmosphere*, 11, 775, <https://doi.org/10.3390/atmos11080775>, 2020.
- 1791 Shuaifeng, S. and Xiaodong, Y.: Projected changes and uncertainty in cold surges over northern China using the CMIP6  
1792 weighted multi-model ensemble, *Atmospheric Res.*, 278, 106334, <https://doi.org/10.1016/j.atmosres.2022.106334>, 2022.
- 1793 Sidhu, B. S., Mehrabi, Z., Ramankutty, N., and Kandlikar, M.: How can machine learning help in understanding the impact  
1794 of climate change on crop yields?, *Environ. Res. Lett.*, 18, 024008, <https://doi.org/10.1088/1748-9326/acb164>, 2023.
- 1795 Sillmann, J., Kharin, V. V., Zhang, X., Zwiers, F. W., and Bronaugh, D.: Climate extremes indices in the CMIP5 multimodel  
1796 ensemble: Part 1. Model evaluation in the present climate, *J. Geophys. Res. Atmospheres*, 118, 1716–1733,  
1797 <https://doi.org/10.1002/jgrd.50203>, 2013.
- 1798 Simpson, I. R., McKinnon, K. A., Davenport, F. V., Tingley, M., Lehner, F., Fahad, A. A., and Chen, D.: Emergent  
1799 Constraints on the Large-Scale Atmospheric Circulation and Regional Hydroclimate: Do They Still Work in CMIP6 and  
1800 How Much Can They Actually Constrain the Future?, *J. Clim.*, 34, 6355–6377, <https://doi.org/10.1175/JCLI-D-21-0055.1>,  
1801 2021.
- 1802 Simpson, I. R., Shaw, T. A., Ceppi, P., Clement, A. C., Fischer, E., Grise, K. M., Pendergrass, A. G., Screen, J. A., Wills, R.  
1803 C. J., Woollings, T., Blackport, R., Kang, J. M., and Po-Chedley, S.: Confronting Earth System Model trends with  
1804 observations, *Sci. Adv.*, 11, eadt8035, <https://doi.org/10.1126/sciadv.adt8035>, 2025.
- 1805 Sippel, S., Kent, E. C., Meinshausen, N., Chan, D., Kadow, C., Neukom, R., Fischer, E. M., Humphrey, V., Rohde, R., De  
1806 Vries, I., and Knutti, R.: Early-twentieth-century cold bias in ocean surface temperature observations, *Nature*, 635, 618–624,  
1807 <https://doi.org/10.1038/s41586-024-08230-1>, 2024.
- 1808 Smith, D. M., Scaife, A. A., Boer, G. J., Caian, M., Doblus-Reyes, F. J., Guemas, V., Hawkins, E., Hazeleger, W.,  
1809 Hermanson, L., Ho, C. K., Ishii, M., Kharin, V., Kimoto, M., Kirtman, B., Lean, J., Matei, D., Merryfield, W. J., Müller, W.  
1810 A., Pohlmann, H., Rosati, A., Wouters, B., and Wyser, K.: Real-time multi-model decadal climate predictions, *Clim. Dyn.*,  
1811 41, 2875–2888, <https://doi.org/10.1007/s00382-012-1600-0>, 2013.
- 1812 Smith, D. M., Eade, R., Andrews, M. B., Ayres, H., Clark, A., Chripko, S., Deser, C., Dunstone, N. J., García-Serrano, J.,  
1813 Gastineau, G., Graff, L. S., Hardiman, S. C., He, B., Hermanson, L., Jung, T., Knight, J., Levine, X., Magnúsdóttir, G.,  
1814 Manzini, E., Matei, D., Mori, M., Msadek, R., Ortega, P., Peings, Y., Scaife, A. A., Screen, J. A., Seabrook, M., Semmler,  
1815 T., Sigmund, M., Streffing, J., Sun, L., and Walsh, A.: Robust but weak winter atmospheric circulation response to future  
1816 Arctic sea ice loss, *Nat. Commun.*, 13, 727, <https://doi.org/10.1038/s41467-022-28283-y>, 2022.

- 1817 Snyder, A., Prime, N., Tebaldi, C., and Dorheim, K.: Uncertainty-informed selection of CMIP6 Earth system model subsets  
 1818 for use in multisectoral and impact models, *Earth Syst. Dyn.*, 15, 1301–1318, <https://doi.org/10.5194/esd-15-1301-2024>,  
 1819 2024.
- 1820 Soares, P. M. M., Careto, J. A. M., Russo, A., and Lima, D. C. A.: The future of Iberian droughts: a deeper analysis based on  
 1821 multi-scenario and a multi-model ensemble approach, *Nat. Hazards*, 117, 2001–2028, [https://doi.org/10.1007/s11069-023-](https://doi.org/10.1007/s11069-023-05938-7)  
 1822 05938-7, 2023.
- 1823 Soares, P. M. M., Johannsen, F., Lima, D. C. A., Lemos, G., Bento, V. A., and Bushenkova, A.: High-resolution  
 1824 downscaling of CMIP6 Earth system and global climate models using deep learning for Iberia, *Geosci. Model Dev.*, 17,  
 1825 229–259, <https://doi.org/10.5194/gmd-17-229-2024>, 2024.
- 1826 Song, X., Wang, D.-Y., Li, F., and Zeng, X.-D.: Evaluating the performance of CMIP6 Earth system models in simulating  
 1827 global vegetation structure and distribution, *Adv. Clim. Change Res.*, 12, 584–595,  
 1828 <https://doi.org/10.1016/j.accre.2021.06.008>, 2021.
- 1829 Sonnewald, M. and Lguensat, R.: Revealing the Impact of Global Heating on North Atlantic Circulation Using Transparent  
 1830 Machine Learning, *J. Adv. Model. Earth Syst.*, 13, e2021MS002496, <https://doi.org/10.1029/2021MS002496>, 2021.
- 1831 Sørland, S. L., Fischer, A. M., Kotlarski, S., Künsch, H. R., Liniger, M. A., Rajczak, J., Schär, C., Spirig, C., Strassmann, K.,  
 1832 and Knutti, R.: CH2018 – National climate scenarios for Switzerland: How to construct consistent multi-model projections  
 1833 from ensembles of opportunity, *Clim. Serv.*, 20, 100196, <https://doi.org/10.1016/j.cliser.2020.100196>, 2020.
- 1834 Steinman, B. A., Frankcombe, L. M., Mann, M. E., Miller, S. K., and England, M. H.: Response to Comment on “Atlantic  
 1835 and Pacific multidecadal oscillations and Northern Hemisphere temperatures,” *Science*, 350, 1326–1326,  
 1836 <https://doi.org/10.1126/science.aac5208>, 2015.
- 1837 Strobach, E. and Bel, G.: Learning algorithms allow for improved reliability and accuracy of global mean surface  
 1838 temperature projections, *Nat. Commun.*, 11, 451, <https://doi.org/10.1038/s41467-020-14342-9>, 2020.
- 1839 Sun, Z. and Archibald, A. T.: Multi-stage ensemble-learning-based model fusion for surface ozone simulations: A focus on  
 1840 CMIP6 models, *Environ. Sci. Ecotechnology*, 8, 100124, <https://doi.org/10.1016/j.ese.2021.100124>, 2021.
- 1841 Tang, B., Hu, W., and Duan, A.: Future Projection of Extreme Precipitation Indices over the Indochina Peninsula and South  
 1842 China in CMIP6 Models, *J. Clim.*, 34, 8793–8811, <https://doi.org/10.1175/JCLI-D-20-0946.1>, 2021.
- 1843 Tang, J., Li, Q., Wang, S., Lee, D.-K., Hui, P., Niu, X., Gutowski, W. J., Dairaku, K., McGregor, J., Katzfey, J., Gao, X.,  
 1844 Wu, J., Hong, S.-Y., Wang, Y., and Sasaki, H.: Building Asian climate change scenario by multi-regional climate models  
 1845 ensemble. Part I: surface air temperature: ASIAN CLIMATE CHANGE BY MULTI-MODEL ENSEMBLE, *Int. J.*  
 1846 *Climatol.*, 36, 4241–4252, <https://doi.org/10.1002/joc.4628>, 2016.
- 1847 Tapiador, F. J., Navarro, A., Moreno, R., Sánchez, J. L., and García-Ortega, E.: Regional climate models: 30 years of  
 1848 dynamical downscaling, *Atmospheric Res.*, 235, 104785, <https://doi.org/10.1016/j.atmosres.2019.104785>, 2020.
- 1849 Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, *J. Geophys. Res. Atmospheres*, 106,  
 1850 7183–7192, <https://doi.org/10.1029/2000JD900719>, 2001.
- 1851 Taylor, M., Caldwell, P. M., Bertagna, L., Clevenger, C., Donahue, A., Foucar, J., Guba, O., Hillman, B., Keen, N., Krishna,  
 1852 J., Norman, M., Sreepathi, S., Terai, C., White, J. B., Salinger, A. G., McCoy, R. B., Leung, L. R., Bader, D. C., and Wu, D.:  
 1853 The Simple Cloud-Resolving E3SM Atmosphere Model Running on the Frontier Exascale System, in: Proceedings of the  
 1854 International Conference for High Performance Computing, Networking, Storage and Analysis, 1–11,  
 1855 <https://doi.org/10.1145/3581784.3627044>, 2023.
- 1856 Tebaldi, C. and Knutti, R.: The use of the multi-model ensemble in probabilistic climate projections, *Philos. Trans. R. Soc.*  
 1857 *Math. Phys. Eng. Sci.*, 365, 2053–2075, <https://doi.org/10.1098/rsta.2007.2076>, 2007.
- 1858 Tebaldi, C., Dorheim, K., Wehner, M., and Leung, R.: Extreme metrics from large ensembles: investigating the effects of  
 1859 ensemble size on their estimates, *Earth Syst. Dyn.*, 12, 1427–1501, <https://doi.org/10.5194/esd-12-1427-2021>, 2021.

- 1860 Tegegne, G., Melesse, A. M., and Worqlul, A. W.: Development of multi-model ensemble approach for enhanced  
1861 assessment of impacts of climate change on climate extremes, *Sci. Total Environ.*, 704, 135357,  
1862 <https://doi.org/10.1016/j.scitotenv.2019.135357>, 2020.
- 1863 Tegegne, G., Melesse, A. M., and Alamirew, T.: Projected changes in extreme precipitation indices from CORDEX  
1864 simulations over Ethiopia, East Africa, *Atmospheric Res.*, 247, 105156, <https://doi.org/10.1016/j.atmosres.2020.105156>,  
1865 2021.
- 1866 Teuling, A. J., de Badts, E. A. G., Jansen, F. A., Fuchs, R., Buitink, J., Hoek van Dijke, A. J., and Sterling, S. M.: Climate  
1867 change, reforestation/afforestation, and urbanization impacts on evapotranspiration and streamflow in Europe, *Hydrol. Earth  
1868 Syst. Sci.*, 23, 3631–3652, <https://doi.org/10.5194/hess-23-3631-2019>, 2019.
- 1869 Thackeray, C. W., Hall, A., Norris, J., and Chen, D.: Constraining the increased frequency of global precipitation extremes  
1870 under warming, *Nat. Clim. Change*, 12, 441–448, <https://doi.org/10.1038/s41558-022-01329-1>, 2022.
- 1871 Thuy, A. and Benoit, D. F.: Explainability through uncertainty: Trustworthy decision-making with neural networks, *Eur. J.  
1872 Oper. Res.*, 317, 330–340, <https://doi.org/10.1016/j.ejor.2023.09.009>, 2024.
- 1873 Toms, B. A., Barnes, E. A., and Ebert-Uphoff, I.: Physically Interpretable Neural Networks for the Geosciences:  
1874 Applications to Earth System Variability, *J. Adv. Model. Earth Syst.*, 12, e2019MS002002,  
1875 <https://doi.org/10.1029/2019MS002002>, 2020.
- 1876 Vázquez-Patiño, A., Campozano, L., Mendoza, D., and Samaniego, E.: A causal flow approach for the evaluation of global  
1877 climate models, *Int. J. Climatol.*, 40, 4497–4517, <https://doi.org/10.1002/joc.6470>, 2020.
- 1878 Veenadhari, S., Misra, B., and Singh, C.: Machine learning approach for forecasting crop yield based on climatic parameters,  
1879 in: 2014 International Conference on Computer Communication and Informatics, 1–5,  
1880 <https://doi.org/10.1109/ICCCI.2014.6921718>, 2014.
- 1881 Vogel, M. M., Hauser, M., and Seneviratne, S. I.: Projected changes in hot, dry and wet extreme events’ clusters in CMIP6  
1882 multi-model ensemble, *Environ. Res. Lett.*, 15, 094021, <https://doi.org/10.1088/1748-9326/ab90a7>, 2020.
- 1883 van Vuuren, D., O’Neill, B., Tebaldi, C., Chini, L., Friedlingstein, P., Hasegawa, T., Riahi, K., Sanderson, B., Govindasamy,  
1884 B., Bauer, N., Eyring, V., Fall, C., Frieler, K., Gidden, M., Gohar, L., Jones, A., King, A., Knutti, R., Kriegler, E., Lawrence,  
1885 P., Lennard, C., Lowe, J., Mathison, C., Mehmood, S., Prado, L., Zhang, Q., Rose, S., Ruane, A., Schleussner, C.-F.,  
1886 Seferian, R., Sillmann, J., Smith, C., Sörensson, A., Panickal, S., Tachiiri, K., Vaughan, N., Vishwanathan, S., Yokohata, T.,  
1887 and Ziehn, T.: The Scenario Model Intercomparison Project for CMIP7 (ScenarioMIP-CMIP7),  
1888 <https://doi.org/10.5194/egusphere-2024-3765>, 30 January 2025.
- 1889 Wang, B., Zheng, L., Liu, D. L., Ji, F., Clark, A., and Yu, Q.: Using multi-model ensembles of CMIP5 global climate models  
1890 to reproduce observed monthly rainfall and temperature with machine learning methods in Australia, *Int. J. Climatol.*, 38,  
1891 4891–4902, <https://doi.org/10.1002/joc.5705>, 2018.
- 1892 Wang, D. and Yuan, F.: High-Performance Computing for Earth System Modeling, in: High Performance Computing for  
1893 Geospatial Applications, edited by: Tang, W. and Wang, S., Springer International Publishing, Cham, 175–184,  
1894 [https://doi.org/10.1007/978-3-030-47998-5\\_10](https://doi.org/10.1007/978-3-030-47998-5_10), 2020.
- 1895 Wang, F. and Tian, D.: On deep learning-based bias correction and downscaling of multiple climate models simulations,  
1896 *Clim. Dyn.*, 59, 3451–3468, <https://doi.org/10.1007/s00382-022-06277-2>, 2022.
- 1897 Wang, F. and Tian, D.: Multivariate bias correction and downscaling of climate models with trend-preserving deep learning,  
1898 *Clim. Dyn.*, 62, 9651–9672, <https://doi.org/10.1007/s00382-024-07406-9>, 2024.
- 1899 Wang, J., Kim, H., Kim, D., Henderson, S. A., Stan, C., and Maloney, E. D.: MJO Teleconnections over the PNA Region in  
1900 Climate Models. Part I: Performance- and Process-Based Skill Metrics, *J. Clim.*, 33, 1051–1067,  
1901 <https://doi.org/10.1175/JCLI-D-19-0253.1>, 2020.
- 1902 Wang, S., Sankaran, S., and Perdikaris, P.: Respecting causality for training physics-informed neural networks, *Comput.*

- 1903 Methods Appl. Mech. Eng., 421, 116813, <https://doi.org/10.1016/j.cma.2024.116813>, 2024.
- 1904 Weber, T., Corotan, A., Hutchinson, B., Kravitz, B., and Link, R.: Technical note: Deep learning for creating surrogate  
1905 models of precipitation in Earth system models, *Atmospheric Chem. Phys.*, 20, 2303–2317, [https://doi.org/10.5194/acp-20-](https://doi.org/10.5194/acp-20-2303-2020)  
1906 2303-2020, 2020.
- 1907 Wehner, M. F.: Characterization of long period return values of extreme daily temperature and precipitation in the CMIP6  
1908 models: Part 2, projections of future change, *Weather Clim. Extrem.*, 30, 100284,  
1909 <https://doi.org/10.1016/j.wace.2020.100284>, 2020.
- 1910 Weigel, A. P., Knutti, R., Liniger, M. A., and Appenzeller, C.: Risks of Model Weighting in Multimodel Climate  
1911 Projections, *J. Clim.*, 23, 4175–4191, <https://doi.org/10.1175/2010JCLI3594.1>, 2010.
- 1912 Wenzel, S., Eyring, V., Gerber, E. P., and Karpechko, A. Yu.: Constraining Future Summer Austral Jet Stream Positions in  
1913 the CMIP5 Ensemble by Process-Oriented Multiple Diagnostic Regression\*, *J. Clim.*, 29, 673–687,  
1914 <https://doi.org/10.1175/JCLI-D-15-0412.1>, 2016.
- 1915 Wilby, R. L. and Fowler, H. J.: *Regional climate downscaling*, Wiley, 85 pp., 2010.
- 1916 Williams, D. N., Balaji, V., Cinquini, L., Denvil, S., Duffy, D., Evans, B., Ferraro, R., Hansen, R., Lautenschlager, M., and  
1917 Trenham, C.: A Global Repository for Planet-Sized Experiments and Observations, *Bull. Am. Meteorol. Soc.*, 97, 803–816,  
1918 <https://doi.org/10.1175/BAMS-D-15-00132.1>, 2016.
- 1919 Wing, A. A., Camargo, S. J., Sobel, A. H., Kim, D., Moon, Y., Murakami, H., Reed, K. A., Vecchi, G. A., Wehner, M. F.,  
1920 Zarzycki, C., and Zhao, M.: Moist Static Energy Budget Analysis of Tropical Cyclone Intensification in High-Resolution  
1921 Climate Models, *J. Clim.*, 32, 6071–6095, <https://doi.org/10.1175/JCLI-D-18-0599.1>, 2019.
- 1922 Woldemeskel, F. M., Sharma, A., Sivakumar, B., and Mehrotra, R.: An error estimation method for precipitation and  
1923 temperature projections for future climates, *J. Geophys. Res. Atmospheres*, 117, <https://doi.org/10.1029/2012JD018062>,  
1924 2012.
- 1925 Wootten, A. M., Başağaoğlu, H., Bertetti, F. P., Chakraborty, D., Sharma, C., Samimi, M., and Mirchi, A.: Customized  
1926 Statistically Downscaled CMIP5 and CMIP6 Projections: Application in the Edwards Aquifer Region in South-Central  
1927 Texas, *Earths Future*, 12, e2024EF004716, <https://doi.org/10.1029/2024EF004716>, 2024.
- 1928 Wu, H., Su, X., and Singh, V. P.: Increasing Risks of Future Compound Climate Extremes With Warming Over Global Land  
1929 Masses, *Earths Future*, 11, e2022EF003466, <https://doi.org/10.1029/2022EF003466>, 2023.
- 1930 Xu, D., Ivanov, V. Y., Kim, J., and Fatichi, S.: On the use of observations in assessment of multi-model climate ensemble,  
1931 *Stoch. Environ. Res. Risk Assess.*, 33, 1923–1937, <https://doi.org/10.1007/s00477-018-1621-2>, 2019.
- 1932 Xu, L. and Wang, A.: Application of the Bias Correction and Spatial Downscaling Algorithm on the Temperature Extremes  
1933 From CMIP5 Multimodel Ensembles in China, *Earth Space Sci.*, 6, 2508–2524, <https://doi.org/10.1029/2019EA000995>,  
1934 2019.
- 1935 Xu, R., Chen, N., Chen, Y., and Chen, Z.: Downscaling and Projection of Multi-CMIP5 Precipitation Using Machine  
1936 Learning Methods in the Upper Han River Basin, *Adv. Meteorol.*, 2020, 8680436, <https://doi.org/10.1155/2020/8680436>,  
1937 2020.
- 1938 Xu, Z., Han, Y., Tam, C.-Y., Yang, Z.-L., and Fu, C.: Bias-corrected CMIP6 global dataset for dynamical downscaling of  
1939 the historical and future climate (1979–2100), *Sci. Data*, 8, 293, <https://doi.org/10.1038/s41597-021-01079-3>, 2021.
- 1940 Yang, T., Hao, X., Shao, Q., Xu, C.-Y., Zhao, C., Chen, X., and Wang, W.: Multi-model ensemble projections in  
1941 temperature and precipitation extremes of the Tibetan Plateau in the 21st century, *Glob. Planet. Change*, 80–81, 1–13,  
1942 <https://doi.org/10.1016/j.gloplacha.2011.08.006>, 2012.
- 1943 Yeganeh-Bakhtiary, A., EyvazOghli, H., Shabakhty, N., Kamranzad, B., and Abolfathi, S.: Machine Learning as a  
1944 Downscaling Approach for Prediction of Wind Characteristics under Future Climate Change Scenarios, *Complexity*, 2022,  
1945 8451812, <https://doi.org/10.1155/2022/8451812>, 2022.

- 1946 Yip, S., Ferro, C. A. T., Stephenson, D. B., and Hawkins, E.: A Simple, Coherent Framework for Partitioning Uncertainty in  
1947 Climate Predictions, *J. Clim.*, 24, 4634–4643, <https://doi.org/10.1175/2011JCLI4085.1>, 2011.
- 1948 Yoon, J. and Schaar, M. van der: E-RNN : Entangled Recurrent Neural Networks for Causal Prediction, 2017.
- 1949 Yu, S., Hannah, W., Peng, L., Lin, J., Bhourri, M. A., Gupta, R., Lütjens, B., Will, J. C., Behrens, G., Busecke, J., Loose, N.,  
1950 Stern, C., Beucler, T., Harrop, B., Hillman, B., Jenney, A., Ferretti, S. L., Liu, N., Anandkumar, A., Brenowitz, N., Eyring,  
1951 V., Geneva, N., Gentine, P., Mandt, S., Pathak, J., Subramaniam, A., Vondrick, C., Yu, R., Zanna, L., Zheng, T.,  
1952 Abernathy, R., Ahmed, F., Bader, D., Baldi, P., Barnes, E., Bretherton, C., Caldwell, P., Chuang, W., Han, Y., Huang, Y.,  
1953 Iglesias-Suarez, F., Jantre, S., Kashinath, K., Khairoutdinov, M., Kurth, T., Lutsko, N., Ma, P.-L., Mooers, G., Neelin, J. D.,  
1954 Randall, D., Shamekh, S., Taylor, M., Urban, N., Yuval, J., Zhang, G., and Pritchard, M.: ClimSim: A large multi-scale  
1955 dataset for hybrid physics-ML climate emulation, *Adv. Neural Inf. Process. Syst.*, 36, 22070–22084, 2023.
- 1956 Zappa, G. and Shepherd, T. G.: Storylines of Atmospheric Circulation Change for European Regional Climate Impact  
1957 Assessment, *J. Clim.*, 30, 6561–6577, <https://doi.org/10.1175/JCLI-D-16-0807.1>, 2017.
- 1958 Zebarjadian, F., Dolatabadi, N., Zahraie, B., Yousefi Sohi, H., and Zandi, O.: Triple coupling random forest approach for  
1959 bias correction of ensemble precipitation data derived from Earth system models for Divandareh-Bijar Basin (Western Iran),  
1960 *Int. J. Climatol.*, 44, 2363–2390, <https://doi.org/10.1002/joc.8458>, 2024.
- 1961 Zhang, X., Zwiers, F. W., Hegerl, G. C., Lambert, F. H., Gillett, N. P., Solomon, S., Stott, P. A., and Nozawa, T.: Detection  
1962 of human influence on twentieth-century precipitation trends, *Nature*, 448, 461–465, <https://doi.org/10.1038/nature06025>,  
1963 2007.
- 1964 Zhang, X., Wang, X.-L., Fan, F., Cheung, Y.-M., and Bose, I.: Enhancing the Performance of Neural Networks Through  
1965 Causal Discovery and Integration of Domain Knowledge, <https://doi.org/10.48550/ARXIV.2311.17303>, 2023.
- 1966 Zhao, L., Wang, Y., Zhao, C., Dong, X., and Yung, Y. L.: Compensating Errors in Cloud Radiative and Physical Properties  
1967 over the Southern Ocean in the CMIP6 Climate Models, *Adv. Atmospheric Sci.*, 39, 2156–2171,  
1968 <https://doi.org/10.1007/s00376-022-2036-z>, 2022.
- 1969 Zhao, T. and Dai, A.: CMIP6 Model-projected Hydroclimatic and Drought Changes and Their Causes in the 21st Century, *J.*  
1970 *Clim.*, 1–58, <https://doi.org/10.1175/JCLI-D-21-0442.1>, 2021.
- 1971 Zhou, W. and Xie, S.-P.: A Hierarchy of Idealized Monsoons in an Intermediate GCM, *J. Clim.*, 31, 9021–9036,  
1972 <https://doi.org/10.1175/JCLI-D-18-0084.1>, 2018.
- 1973 Zhu, J. and Poulsen, C. J.: Last Glacial Maximum (LGM) climate forcing and ocean dynamical feedback and their  
1974 implications for estimating climate sensitivity, *Clim. Past*, 17, 253–267, <https://doi.org/10.5194/cp-17-253-2021>, 2021.
- 1975 Zuluaga, M., Sergent, G., Krause, A., and Püschel, M.: Active Learning for Multi-Objective Optimization, in: Proceedings of  
1976 the 30th International Conference on Machine Learning, 462–470, 2013.

## 1977 **Appendix**

### 1978 **A Statistics of the field over past decades**

1979 Figures 5 and 6 were built using data from the Web of Science database. The queries for each category are:

#### 1980 **Total ML:**

1981 TS=("machine learning" OR "artificial intelligence" OR "neural networks" OR "random forest" OR "decision trees" OR  
1982 "deep learning" OR "supervised learning" OR "unsupervised learning") AND TS=("CMIP" OR "CMIP3" OR "CMIP5" OR

1983 "CMIP6" OR "Coupled Model Intercomparison Project" OR "climate model" OR "climate models" OR "general circulation  
1984 model" OR "general circulation models" OR "Earth system model" OR "Earth system models")

1985 **ML-MME:**

1986 TS=("machine learning" OR "artificial intelligence" OR "neural networks" OR "random forest" OR "decision trees" OR  
1987 "deep learning" OR "supervised learning" OR "unsupervised learning") AND TS=("CMIP" OR "CMIP3" OR "CMIP5" OR  
1988 "CMIP6" OR "Coupled Model Intercomparison Project" OR "climate model" OR "climate models" OR "general circulation  
1989 model" OR "general circulation models" OR "Earth system model" OR "Earth system models") AND TS=("multi-model  
1990 ensemble" OR " multi-model ensembles")

1991 **ML-Downscaling:**

1992 TS=("machine learning" OR "artificial intelligence" OR "neural networks" OR "random forest" OR "decision trees" OR  
1993 "deep learning" OR "supervised learning" OR "unsupervised learning") AND TS=("CMIP" OR "CMIP3" OR "CMIP5" OR  
1994 "CMIP6" OR "Coupled Model Intercomparison Project" OR "climate model" OR "climate models" OR "general circulation  
1995 model" OR "general circulation models" OR "Earth system model" OR "Earth system models") AND TS=("downscaling"  
1996 OR "bias correction")

1997 **ML-Downscaling MME:**

1998 TS=("machine learning" OR "artificial intelligence" OR "neural networks" OR "random forest" OR "decision trees" OR  
1999 "deep learning" OR "supervised learning" OR "unsupervised learning") AND TS=("CMIP" OR "CMIP3" OR "CMIP5" OR  
2000 "CMIP6" OR "Coupled Model Intercomparison Project" OR "climate model" OR "climate models" OR "general circulation  
2001 model" OR "general circulation models" OR "Earth system model" OR "Earth system models") AND TS=("downscaling"  
2002 OR "bias correction") AND TS=("multi-model ensemble" OR " multi-model ensembles")

2003 **ML Causality:**

2004 TS=("CMIP" OR "CMIP3" OR "CMIP5" OR "CMIP6" OR "Coupled Model Intercomparison Project" OR "climate model"  
2005 OR "climate models" OR "general circulation model" OR "general circulation models" OR "Earth system model" OR "Earth  
2006 system models") AND TS=("causal discovery" OR "causality" OR "causal inference" OR "causal")

2007 **ML Emulators:**

2008 TS=("machine learning" OR "artificial intelligence" OR "neural networks" OR "random forest" OR "decision trees" OR  
2009 "deep learning" OR "supervised learning" OR "unsupervised learning") AND TS=("CMIP" OR "CMIP3" OR "CMIP5" OR  
2010 "CMIP6" OR "Coupled Model Intercomparison Project" OR "climate model" OR "climate models" OR "general circulation  
2011 model" OR "general circulation models" OR "Earth system model" OR "Earth system models") AND TS=("emulation" or  
2012 "surrogate" or "emulator" or "emulators" or "surrogates")

### 2013 **ML XAI:**

2014 TS=("machine learning" OR "artificial intelligence" OR "neural networks" OR "random forest" OR "decision trees" OR  
2015 "deep learning" OR "supervised learning" OR "unsupervised learning") AND TS=("CMIP" OR "CMIP3" OR "CMIP5" OR  
2016 "CMIP6" OR "Coupled Model Intercomparison Project" OR "climate model" OR "climate models" OR "general circulation  
2017 model" OR "general circulation models" OR "Earth system model" OR "Earth system models") AND TS=( "XAI" OR  
2018 "explainable AI" OR "Layer-wise Relevance Propagation" OR "LRP" OR "Feature importance analysis" OR "feature  
2019 importance")

### 2020 **Model Independence:**

2021 TS=("climate" OR "Earth" OR "Earth System") AND TS=("CMIP" OR "Coupled Model Intercomparison Project" OR  
2022 "climate model" OR "general circulation model") AND TS=("ensemble" OR "multi-model ensemble") AND  
2023 TS=("dependence" OR "independence" OR "genealogy")

### 2024 **SMILEs:**

2025 TS=("Multi-model ensemble" OR "coupled model intercomparison" OR "cmip") AND TS= ("large ensemble" OR "grand  
2026 ensemble" OR "smile")

### 2027 **B Systematic Model Biases**

2028 Some systematic biases are present in the vast majority of CMIP models at the global and regional scale, and some might even  
2029 persist over multiple CMIP generations, which requires special attention. In this section we review some long-standing biases  
2030 in CMIP models and strive to discuss their origins and consequences of these systematic model biases. Those are as follows:  
2031 a) General evaluation, b) Sea surface temperature (SST) and ocean model biases, c) The Intertropical Convergence Zone  
2032 (ITCZ) bias, d) biases in extratropical cyclones, e) Marine tropical/subtropical low cloud biases, f) biases in the cryosphere, g)  
2033 biases in extremes. With this list, we do not intend to provide a complete list of all bias reported, but to give some relevant  
2034 examples of model biases and its background. For further details on this topic, we also recommend Simpson et al. (2025).

2035 *General evaluation:* Bock et al., 2020 employed the ESMValTool (see Section 2.6 and Eyring et al., 2020; Righi et al., 2020),  
2036 to quantify the progress of climate models across different CMIP phases. Their analysis revealed significant advancements  
2037 from CMIP3 to CMIP6 in simulating the vertical distributions of key variables, including temperature, water vapor, and zonal  
2038 wind speed. The authors also demonstrated that high-resolution models in the historical CMIP6 simulations show a notable  
2039 reduction of temperature and precipitation mean biases.

2040 *Sea surface temperature and ocean model biases:* The ocean accumulates more than 90% of the excess energy from the global  
2041 greenhouse effect (IPCC, AR6). The oceanic global circulation gyres transport excess heat from the tropics towards the poles.  
2042 Furthermore, the oceanic surface fluxes of heat and moisture enter the atmosphere and thereby affect its dynamics. The ocean  
2043 component also interacts with the cryosphere and influences processes therein (IPCC, AR6). These various oceanic processes  
2044 have to be properly captured in ESMs. Long-standing SST biases result in biases when simulating other key phenomena such  
2045 as tropical cyclones (e.g. Duteuil et al., 2020) and extratropical cyclones (e.g., Priestley et al., 2023a). Wills et al. (2022)  
2046 investigated systematic biases in the large-scale patterns of recent SST and sea-level pressure change and showed that CMIP5  
2047 and CMIP6 ensembles are not able to reproduce the observed trends. Luo et al. (2023), moreover, discussed the origins of  
2048 Southern Ocean warm SST bias in CMIP6 models. The Southern Ocean has namely been subjected to systematic warm SST  
2049 bias in several generations of CMIP models (Sen Gupta et al., 2009; Wang et al., 2014). Westen and Dijkstra (2024) recently  
2050 discussed persistent climate model biases in the Atlantic Ocean's freshwater transport. These various aforementioned biases  
2051 are linked to the Atlantic Meridional Overturning Circulation (AMOC), which consists of the northward flow in the upper  
2052 oceanic layers and returning southward flow in the deep ocean (Luo et al., 2023; Wang et al., 2024). The AMOC is considered  
2053 to be one of the major tipping elements in the global climate system (Armstrong McKay et al., 2022; Van Westen et al., 2024),  
2054 which may weaken or even collapse with future global warming, thus a more reliable representation of SST/ocean model  
2055 would be desirable e.g. to better foresee the future AMOC behaviour.

2056 *The Intertropical Convergence Zone (ITCZ) bias:* ITCZ is a band of a zonally-oriented surface convergence zone near the  
2057 equator associated with deep convective clouds and heavy precipitation (Schneider et al., 2014; Waliser and Gautier, 1993).  
2058 The common problem of fully-coupled global climate models from the early stage of their development is that they simulate  
2059 two ITCZs over the central and eastern Pacific and the Atlantic in both hemispheres, instead of one ITCZ over the northern  
2060 hemisphere as in observations, which is referred to as the double-ITCZ bias (Adam et al., 2018; Li and Xie, 2014; Oueslati  
2061 and Bellon, 2015; Tian and Dong, 2020; Xiang et al., 2017). Tian and Dong (2020), as an illustration, recently examined the  
2062 double-ITCZ bias in CMIP3, CMIP5, and CMIP6 based on annual mean precipitation. They found that all three generations  
2063 of CMIP models exhibit similar systematic annual MME mean precipitation errors in the tropics when evaluated against the  
2064 NOAA Global Precipitation Climatology Project (GPCP; Adler et al., 2003) and the NASA Tropical Rainfall Measurement  
2065 Mission (TRMM; Huffman et al., 2007) observational datasets.

2066 *Biases in extratropical cyclones:* Extratropical cyclones involving weather fronts and related overall storm tracks are an  
2067 important component of the climate system since they transport heat poleward and are associated with a notable amount of  
2068 precipitation and severe weather in the midlatitudes (Clark and Gray, 2020; Dacre, 2020; Schultz et al., 2019). The accurate  
2069 representation of extratropical cyclones, including their thermodynamics, frontal structure, and track in CMIP models,  
2070 however, remains challenging and has been subjected to biases (e.g. Chang et al., 2012; Priestley et al., 2023a, b). Priestley et  
2071 al. (2023a) investigated drivers of biases in the CMIP6 extratropical storm tracks in the Northern Hemisphere (NH). Even  
2072 though the previous work demonstrated that the representation of extratropical storm tracks in the NH has improved from  
2073 CMIP5 to CMIP6, the persistent biases remain in CMIP6 (Priestley et al., 2023a). A follow-up study by Priestley et al. (2023b)  
2074 investigated drivers of biases in the CMIP6 extratropical storm tracks in the Southern Hemisphere (SH). The Southern  
2075 Hemisphere storm tracks have been commonly simulated too far equatorward in CMIP models during the historical period.  
2076 This issue was somewhat reduced in CMIP6 compared to CMIP5, although it is still a problem.

2077 *Marine tropical/subtropical low cloud biases:* Črnivec et al. (2023) analyzed 12 CMIP6 ESMs and demonstrated that they all  
2078 underestimate the aerial extent of low clouds and simultaneously overestimate their radiative effect at the top of the atmosphere.  
2079 This well-known issue, referred to as the “too few, too bright” tropical low-cloud bias, was already present in previous  
2080 generations of climate models such as CMIP5 and CMIP3 (e.g., Nam et al., 2012, and references therein). Cesana et al. (2023),  
2081 moreover, addressed how the representation of marine tropical Sc and Cu clouds and associated feedbacks in the abrupt 4xCO2  
2082 scenario changed between CMIP5 and CMIP6. They found that, collectively, CMIP6 models notably increased Sc cloud cover  
2083 and slightly increased Cu cloud cover compared to their CMIP5 predecessors and are thus closer to observations. They further  
2084 showed that CMIP6 models notably improved the representation of Sc feedback and slightly improved the representation of  
2085 Cu feedback compared to CMIP5 models. Yet CMIP6 models still underestimate the magnitude of positive Sc and Cu  
2086 feedbacks relative to observationally inferred estimates, which should drive further climate model development.

2087 *Biases in the cryosphere:* The global cryosphere plays an important role in determining the planetary climate since bright ice  
2088 and snow surfaces reflect a significant portion of the solar radiation back to space and cool the planet (IPCC, AR6). In a  
2089 warming world, sea ice is shrinking and thinning, with both Arctic and Antarctic sea ice approaching historic lows (NASA  
2090 Earth Observatory; IPCC AR6). The melting of sea ice with global surface warming implies that an increasing area of dark  
2091 and absorptive ocean surface is exposed to warming sunlight, which forms one of the principal climate feedback mechanisms  
2092 – namely, the sea ice albedo feedback (IPCC, AR6). It is thus pivotal to best capture the cryosphere extent, properties, and its  
2093 response to global warming. To that end, Frankignoul et al. (2024) investigated Arctic September sea ice concentration biases  
2094 in CMIP6 models and their relationships with other model variables. They demonstrated that CMIP6 models exhibit large  
2095 biases in Arctic sea ice climatology, which seem to be related to biases in seasonal oceanic and atmospheric circulations. Notz  
2096 and the Sea-Ice Model Intercomparison Project (SIMIP) Community (2020) furthermore showed that CMIP6 models still fail  
2097 to simulate a plausible evolution of Arctic sea-ice area (SIA), even though CMIP6 models better capture the sensitivity of

2098 Arctic sea ice to forcing changes compared to CMIP5 and CMIP3 models. Roach et al. (2020) evaluated the Antarctic sea ice  
2099 in CMIP6 and demonstrated that the mean Antarctic sea-ice area is close to satellite observations, but inter-model spread  
2100 remains substantial, with summer Antarctic SIA being consistently biased low across the ensemble. Nevertheless, they found  
2101 modest improvements in the simulation of sea-ice area and concentration compared to CMIP5.

## 2102 **C Further Examples of process-based Evaluation**

2103 Another example are low-level clouds over tropical and subtropical oceans that have been poorly simulated in multiple CMIP  
2104 generations when evaluated against satellite observations in the present-day climate (e.g. Nam et al., 2012), which inhibits  
2105 reliable future climate projections. Črnivec et al. (2023) and Cesana et al. (2023) introduced a qualitative approach to  
2106 discriminate stratocumulus (Sc) from shallow cumulus (Cu) low-cloud regimes to evaluate their horizontal extent (cloud  
2107 cover), radiative effect at the top of the atmosphere (TOA) and cloud-radiative feedbacks in CMIP5 and CMIP6 models. This  
2108 approach is essential for guiding model improvements, because Sc and Cu formation and evolution are driven by a distinct  
2109 interplay of coupled processes within the moist marine boundary layer (such as radiation, turbulence, convection); and Sc and  
2110 Cu clouds also respond differently to global warming (Cesana and Del Genio, 2021).

2111 The teleconnection between the Indian Summer Monsoon (ISM) and the El Niño–Southern Oscillation (ENSO) serves as a  
2112 further example, which is well captured by MMEs of CMIP5 and CMIP6 models (Roy and Tedeschi, 2016; Roy et al., 2017).  
2113 The teleconnection is strongest over central northeast India (Roy et al., 2017), where El Niño events are associated with a  
2114 significant rainfall deficit, while La Niña events lead to a significant rainfall excess. In these studies, the MME is constructed  
2115 using a simple mean (“one-model-one-vote”) approach. Similar results are obtained when the MME is restricted to a subset of  
2116 well-performing models, as identified for the ISM by Jourdain et al. (2013). Precipitation anomalies associated with different  
2117 ENSO phases are reproduced well by most individual models and by the MME, in agreement with observations (see Roy et  
2118 al., 2017 for details). Furthermore, the model ensemble analysis of ISM precipitation and Pacific SSTs reveals a clear linkage  
2119 between the Walker circulation and the ISM over central northeast India, consistent with observations. This region represents  
2120 the meeting point of the Hadley and Walker circulations during the ISM season, and the associated coupling and teleconnection  
2121 processes appear to be well represented in most CMIP models and in the MME. This process-based understanding helps  
2122 explain why the ENSO–ISM teleconnection is robustly captured in climate models.

2123 Ahmed and Neelin (2021) utilised the observed relationship between tropical precipitation and buoyancy as the basis for a  
2124 process-oriented analysis of CMIP6 models. They quantitatively assessed the thermodynamic sensitivity of convection across  
2125 models using regime-oriented diagnostics. Their results showed that several models exhibit excessive moisture sensitivity,  
2126 potentially arising from underactive convective schemes or tuning assumptions. Consequently, these models tend to produce  
2127 mean precipitation states that are biased towards grid-scale saturation.

## 2128 **D Storylines for understanding Uncertainty**

2129 Considering this, uncertainty in climate projections can be communicated through climate storylines (Shepherd et al., 2018),  
2130 which emphasizes exploring and understanding physically plausible events or pathways. The storyline approach differs from  
2131 traditional methods of uncertainty evaluation in climate models in that it does not assume that model spread adequately  
2132 represents uncertainty. This assumption may not hold for dynamically driven climate phenomena, where MME means may  
2133 obscure regional details with individual climate models exhibiting atmospheric circulation patterns that can differ qualitatively  
2134 from the multi-model mean (Bellomo et al., 2021; Zappa and Shepherd, 2017). Instead of quantifying the likelihood of events,  
2135 storylines focus on the physical drivers and interactions that make an event possible (Shepherd et al., 2018), constructing a  
2136 causal network and conditioning on specific assumptions. If we know thermodynamic changes are robust, the thermodynamic  
2137 aspects of the observed changes are regarded as certain and the dynamic aspects as uncertain. By explicitly linking causal  
2138 mechanisms to regional climate hazards, storylines are especially useful for regional climate impacts and understanding  
2139 extreme events (Bevacqua et al., 2022; Shepherd, 2019; Zappa and Shepherd, 2017), improving the interpretability and  
2140 usability of projections for decision-makers (Kunimitsu et al., 2023).

## 2141 **References**

- 2142 Adam, O., Schneider, T., and Brient, F.: Regional and seasonal variations of the double-ITCZ bias in CMIP5 models, *Clim*  
2143 *Dyn*, 51, 101–117, <https://doi.org/10.1007/s00382-017-3909-1>, 2018.
- 2144 Adler, R. F., Huffman, G. J., Chang, A., Ferraro, R., Xie, P.-P., Janowiak, J., Rudolf, B., Schneider, U., Curtis, S., Bolvin, D.,  
2145 Gruber, A., Susskind, J., Arkin, P., and Nelkin, E.: The Version-2 Global Precipitation Climatology Project (GPCP) Monthly  
2146 Precipitation Analysis (1979–Present), *J. Hydrometeor*, 4, 1147–1167, [https://doi.org/10.1175/1525-7541\(2003\)004%253C1147:TVGPCP%253E2.0.CO;2](https://doi.org/10.1175/1525-7541(2003)004%253C1147:TVGPCP%253E2.0.CO;2), 2003.
- 2148 Ahmed, F. and Neelin, J. D.: A Process-Oriented Diagnostic to Assess Precipitation-Thermodynamic Relations and  
2149 Application to CMIP6 Models, *Geophysical Research Letters*, 48, e2021GL094108, <https://doi.org/10.1029/2021GL094108>,  
2150 2021.
- 2151 Armstrong McKay, D. I., Staal, A., Abrams, J. F., Winkelmann, R., Sakschewski, B., Loriani, S., Fetzer, I., Cornell, S. E.,  
2152 Rockström, J., and Lenton, T. M.: Exceeding 1.5°C global warming could trigger multiple climate tipping points, *Science*,  
2153 377, eabn7950, <https://doi.org/10.1126/science.abn7950>, 2022.
- 2154 Bellomo, K., Angeloni, M., Corti, S., and von Hardenberg, J.: Future climate change shaped by inter-model differences in  
2155 Atlantic meridional overturning circulation response, *Nat Commun*, 12, 3659, <https://doi.org/10.1038/s41467-021-24015-w>,  
2156 2021.
- 2157 Bevacqua, E., Zappa, G., Lehner, F., and Zscheischler, J.: Precipitation trends determine future occurrences of compound hot-  
2158 dry events, *Nat. Clim. Chang.*, 12, 350–355, <https://doi.org/10.1038/s41558-022-01309-5>, 2022.
- 2159 Bock, L., Lauer, A., Schlund, M., Barreiro, M., Bellouin, N., Jones, C., Meehl, G. A., Predoi, V., Roberts, M. J., and Eyring,  
2160 V.: Quantifying Progress Across Different CMIP Phases With the ESMValTool, *JGR Atmospheres*, 125, e2019JD032321,  
2161 <https://doi.org/10.1029/2019JD032321>, 2020.
- 2162 Cesana, G. V. and Del Genio, A. D.: Observational constraint on cloud feedbacks suggests moderate climate sensitivity, *Nat.*  
2163 *Clim. Chang.*, 11, 213–218, <https://doi.org/10.1038/s41558-020-00970-y>, 2021.

- 2164 Cesana, G. V., Ackerman, A. S., Črnivec, N., Pincus, R., and Chepfer, H.: An observation-based method to assess tropical  
2165 stratocumulus and shallow cumulus clouds and feedbacks in CMIP6 and CMIP5 models, *Environ. Res. Commun.*, 5, 045001,  
2166 <https://doi.org/10.1088/2515-7620/acc78a>, 2023.
- 2167 Chang, E. K. M., Guo, Y., and Xia, X.: CMIP5 multimodel ensemble projection of storm track change under global warming,  
2168 *J. Geophys. Res.*, 117, 2012JD018578, <https://doi.org/10.1029/2012JD018578>, 2012.
- 2169 Clark, P. A. and Gray, S. L.: Sting jets in extratropical cyclones: a review, *Quart J Royal Meteorol Soc*, 146, 1065–1065,  
2170 <https://doi.org/10.1002/qj.3761>, 2020.
- 2171 Črnivec, N., Cesana, G., and Pincus, R.: Evaluating the Representation of Tropical Stratocumulus and Shallow Cumulus  
2172 Clouds As Well As Their Radiative Effects in CMIP6 Models Using Satellite Observations, *JGR Atmospheres*, 128,  
2173 e2022JD038437, <https://doi.org/10.1029/2022JD038437>, 2023.
- 2174 Dacre, H. F.: A review of extratropical cyclones: observations and conceptual models over the past 100 years, *Weather*, 75,  
2175 4–7, <https://doi.org/10.1002/wea.3653>, 2020.
- 2176 Dutheil, C., Lengaigne, M., Bador, M., Vialard, J., Lefèvre, J., Jourdain, N. C., Jullien, S., Peltier, A., Sultan, B., and Menkès,  
2177 C.: Impact of projected sea surface temperature biases on tropical cyclones projections in the South Pacific, *Sci Rep*, 10, 4838,  
2178 <https://doi.org/10.1038/s41598-020-61570-6>, 2020.
- 2179 Eyring, V., Bock, L., Lauer, A., Righi, M., Schlund, M., Andela, B., Arnone, E., Bellprat, O., Brötz, B., Caron, L.-P.,  
2180 Carvalhais, N., Cionni, I., Cortesi, N., Crezee, B., Davin, E. L., Davini, P., Debeire, K., De Mora, L., Deser, C., Docquier, D.,  
2181 Earnshaw, P., Ehbrecht, C., Gier, B. K., Gonzalez-Reviriego, N., Goodman, P., Hagemann, S., Hardiman, S., Hassler, B.,  
2182 Hunter, A., Kadow, C., Kindermann, S., Koirala, S., Koldunov, N., Lejeune, Q., Lembo, V., Lovato, T., Lucarini, V.,  
2183 Massonnet, F., Müller, B., Pandde, A., Pérez-Zanón, N., Phillips, A., Predoi, V., Russell, J., Sellar, A., Serva, F., Stacke, T.,  
2184 Swaminathan, R., Torralba, V., Vegas-Regidor, J., Von Hardenberg, J., Weigel, K., and Zimmermann, K.: Earth System Model  
2185 Evaluation Tool (ESMValTool) v2.0 – an extended set of large-scale diagnostics for quasi-operational and comprehensive  
2186 evaluation of Earth system models in CMIP, *Geosci. Model Dev.*, 13, 3383–3438, <https://doi.org/10.5194/gmd-13-3383-2020>,  
2187 2020.
- 2188 Frankignoul, C., Raillard, L., Ferster, B., and Kwon, Y.-O.: Arctic September Sea Ice Concentration Biases in CMIP6 Models  
2189 and Their Relationships with Other Model Variables, *Journal of Climate*, 37, 4257–4274, <https://doi.org/10.1175/JCLI-D-23-0452.1>, 2024.
- 2191 Huffman, G. J., Bolvin, D. T., Nelkin, E. J., Wolff, D. B., Adler, R. F., Gu, G., Hong, Y., Bowman, K. P., and Stocker, E. F.:  
2192 The TRMM Multisatellite Precipitation Analysis (TMPA): Quasi-Global, Multiyear, Combined-Sensor Precipitation  
2193 Estimates at Fine Scales, *Journal of Hydrometeorology*, 8, 38–55, <https://doi.org/10.1175/JHM560.1>, 2007.
- 2194 Intergovernmental Panel on Climate Change (IPCC): Climate Change 2021 – The Physical Science Basis: Working Group I  
2195 Contribution to the Sixth Assessment Report of the Intergovernmental Panel on Climate Change, Cambridge University Press,  
2196 Cambridge, <https://doi.org/10.1017/9781009157896>, 2021.
- 2197 Jourdain, N. C., Gupta, A. S., Taschetto, A. S., Ummenhofer, C. C., Moise, A. F., and Ashok, K.: The Indo-Australian monsoon  
2198 and its relationship to ENSO and IOD in reanalysis data and the CMIP3/CMIP5 simulations, *Clim Dyn*, 41, 3073–3102,  
2199 <https://doi.org/10.1007/s00382-013-1676-1>, 2013.
- 2200 Kunimitsu, T., Baldissera Pacchetti, M., Ciullo, A., Sillmann, J., Shepherd, T. G., Taner, M. Ü., and van den Hurk, B.:  
2201 Representing storylines with causal networks to support decision making: Framework and example, *Climate Risk  
2202 Management*, 40, 100496, <https://doi.org/10.1016/j.crm.2023.100496>, 2023.
- 2203 Li, G. and Xie, S.-P.: Tropical Biases in CMIP5 Multimodel Ensemble: The Excessive Equatorial Pacific Cold Tongue and  
2204 Double ITCZ Problems\*, *Journal of Climate*, 27, 1765–1780, <https://doi.org/10.1175/JCLI-D-13-00337.1>, 2014.
- 2205 Luo, F., Ying, J., Liu, T., and Chen, D.: Origins of Southern Ocean warm sea surface temperature bias in CMIP6 models, *npj  
2206 Clim Atmos Sci*, 6, 1–8, <https://doi.org/10.1038/s41612-023-00456-6>, 2023.

- 2207 Nam, C., Bony, S., Dufresne, J. -L., and Chepfer, H.: The ‘too few, too bright’ tropical low-cloud problem in CMIP5 models,  
2208 *Geophysical Research Letters*, 39, 2012GL053421, <https://doi.org/10.1029/2012GL053421>, 2012.
- 2209 Oueslati, B. and Bellon, G.: The double ITCZ bias in CMIP5 models: interaction between SST, large-scale circulation and  
2210 precipitation, *Climate Dynamics*, 44, 585–607, <https://doi.org/10.1007/s00382-015-2468-6>, 2015.
- 2211 Priestley, M. D. K., Ackerley, D., Catto, J. L., and Hodges, K. I.: Drivers of Biases in the CMIP6 Extratropical Storm Tracks.  
2212 Part I: Northern Hemisphere, *Journal of Climate*, 36, 1451–1467, <https://doi.org/10.1175/JCLI-D-20-0976.1>, 2023a.
- 2213 Priestley, M. D. K., Ackerley, D., Catto, J. L., and Hodges, K. I.: Drivers of Biases in the CMIP6 Extratropical Storm Tracks.  
2214 Part II: Southern Hemisphere, *Journal of Climate*, 36, 1469–1486, <https://doi.org/10.1175/JCLI-D-20-0977.1>, 2023b.
- 2215 Righi, M., Andela, B., Eyring, V., Lauer, A., Predoi, V., Schlund, M., Vegas-Regidor, J., Bock, L., Brötz, B., De Mora, L.,  
2216 Diblen, F., Dreyer, L., Drost, N., Earnshaw, P., Hassler, B., Koldunov, N., Little, B., Loosveldt Tomas, S., and Zimmermann,  
2217 K.: Earth System Model Evaluation Tool (ESMValTool) v2.0 – technical overview, *Geosci. Model Dev.*, 13, 1179–1199,  
2218 <https://doi.org/10.5194/gmd-13-1179-2020>, 2020.
- 2219 Roach, L. A., Dörr, J., Holmes, C. R., Massonnet, F., Blockley, E. W., Notz, D., Rackow, T., Raphael, M. N., O’Farrell, S. P.,  
2220 Bailey, D. A., and Bitz, C. M.: Antarctic Sea Ice Area in CMIP6, *Geophysical Research Letters*, 47, e2019GL086729,  
2221 <https://doi.org/10.1029/2019GL086729>, 2020.
- 2222 Roy, I. and Tedeschi, R.: Influence of ENSO on Regional Indian Summer Monsoon Precipitation—Local Atmospheric  
2223 Influences or Remote Influence from Pacific, *Atmosphere*, 7, 25, <https://doi.org/10.3390/atmos7020025>, 2016.
- 2224 Roy, I., Tedeschi, R. G., and Collins, M.: ENSO teleconnections to the Indian summer monsoon in observations and models,  
2225 *Intl Journal of Climatology*, 37, 1794–1813, <https://doi.org/10.1002/joc.4811>, 2017.
- 2226 Schneider, T., Bischoff, T., and Haug, G. H.: Migrations and dynamics of the intertropical convergence zone, *Nature*, 513, 45–  
2227 53, <https://doi.org/10.1038/nature13636>, 2014.
- 2228 Schultz, D. M., Bosart, L. F., Colle, B. A., Davies, H. C., Dearden, C., Keyser, D., Martius, O., Roebber, P. J., Steenburgh, W.  
2229 J., Volkert, H., and Winters, A. C.: Extratropical Cyclones: A Century of Research on Meteorology’s Centerpiece,  
2230 *Meteorological Monographs*, 59, 16.1-16.56, <https://doi.org/10.1175/AMSMONOGRAPHS-D-18-0015.1>, 2019.
- 2231 Sen Gupta, A., Santoso, A., Taschetto, A. S., Ummenhofer, C. C., Trevena, J., and England, M. H.: Projected Changes to the  
2232 Southern Hemisphere Ocean and Sea Ice in the IPCC AR4 Climate Models, *Journal of Climate*, 22, 3047–3078,  
2233 <https://doi.org/10.1175/2008JCLI2827.1>, 2009.
- 2234 Shepherd, T. G.: Storyline approach to the construction of regional climate change information, *Proc. R. Soc. A.*, 475,  
2235 20190013, <https://doi.org/10.1098/rspa.2019.0013>, 2019.
- 2236 Shepherd, T. G., Boyd, E., Calel, R. A., Chapman, S. C., Dessai, S., Dima-West, I. M., Fowler, H. J., James, R., Maraun, D.,  
2237 Martius, O., Senior, C. A., Sobel, A. H., Stainforth, D. A., Tett, S. F. B., Trenberth, K. E., Van Den Hurk, B. J. J. M., Watkins,  
2238 N. W., Wilby, R. L., and Zenghelis, D. A.: Storylines: an alternative approach to representing uncertainty in physical aspects  
2239 of climate change, *Climatic Change*, 151, 555–571, <https://doi.org/10.1007/s10584-018-2317-9>, 2018.
- 2240 Simpson, I. R., Shaw, T. A., Ceppi, P., Clement, A. C., Fischer, E., Grise, K. M., Pendergrass, A. G., Screen, J. A., Wills, R.  
2241 C. J., Woollings, T., Blackport, R., Kang, J. M., and Po-Chedley, S.: Confronting Earth System Model trends with observations,  
2242 *Sci. Adv.*, 11, eadt8035, <https://doi.org/10.1126/sciadv.adt8035>, 2025.
- 2243 Tian, B. and Dong, X.: The Double-ITCZ Bias in CMIP3, CMIP5, and CMIP6 Models Based on Annual Mean Precipitation,  
2244 *Geophysical Research Letters*, 47, e2020GL087232, <https://doi.org/10.1029/2020GL087232>, 2020.
- 2245 Van Westen, R. M., Kliphuis, M., and Dijkstra, H. A.: Physics-based early warning signal shows that AMOC is on tipping  
2246 course, *Sci. Adv.*, 10, eadk1189, <https://doi.org/10.1126/sciadv.adk1189>, 2024.
- 2247 Waliser, D. E. and Gautier, C.: A Satellite-derived Climatology of the ITCZ, *Journal of Climate*, 6, 2162–2174,  
2248 [https://doi.org/10.1175/1520-0442\(1993\)006%253C2162:ASDCOT%253E2.0.CO;2](https://doi.org/10.1175/1520-0442(1993)006%253C2162:ASDCOT%253E2.0.CO;2), 1993.

- 2249 Wang, C., Zhang, L., Lee, S.-K., Wu, L., and Mechoso, C. R.: A global perspective on CMIP5 climate model biases, *Nature*  
2250 *Clim Change*, 4, 201–205, <https://doi.org/10.1038/nclimate2118>, 2014.
- 2251 Wang, S., Sankaran, S., and Perdikaris, P.: Respecting causality for training physics-informed neural networks, *Computer*  
2252 *Methods in Applied Mechanics and Engineering*, 421, 116813, <https://doi.org/10.1016/j.cma.2024.116813>, 2024.
- 2253 van Westen, R. M. and Dijkstra, H. A.: Persistent climate model biases in the Atlantic Ocean’s freshwater transport, *Ocean*  
2254 *Science*, 20, 549–567, <https://doi.org/10.5194/os-20-549-2024>, 2024.
- 2255 Wills, R. C. J., Dong, Y., Proistosescu, C., Armour, K. C., and Battisti, D. S.: Systematic Climate Model Biases in the Large-  
2256 Scale Patterns of Recent Sea-Surface Temperature and Sea-Level Pressure Change, *Geophysical Research Letters*, 49,  
2257 e2022GL100011, <https://doi.org/10.1029/2022GL100011>, 2022.
- 2258 Xiang, B., Zhao, M., Held, I. M., and Golaz, J.: Predicting the severity of spurious “double ITCZ” problem in CMIP5 coupled  
2259 models from AMIP simulations, *Geophysical Research Letters*, 44, 1520–1527, <https://doi.org/10.1002/2016GL071992>, 2017.
- 2260 Zappa, G. and Shepherd, T. G.: Storylines of Atmospheric Circulation Change for European Regional Climate Impact  
2261 Assessment, *Journal of Climate*, 30, 6561–6577, <https://doi.org/10.1175/JCLI-D-16-0807.1>, 2017.