

1 **Systematic and Objective Evaluation of Earth System Models: PCMDI**
2 **Metrics Package (PMP) version 3**

3

4 Jiwoo Lee¹, Peter J. Gleckler¹, Min-Seop Ahn^{2,3}, Ana Ordonez¹, Paul A. Ullrich^{1,4}, Kenneth R.
5 Sperber^{1,a}, Karl E. Taylor¹, Yann Y. Planton^{5,6}, Eric Guilyardi^{7,8}, Paul Durack¹, Celine Bonfils¹,
6 Mark D. Zelinka¹, Li-Wei Chao¹, Bo Dong¹, Charles Doutriaux¹, Chengzhu Zhang¹, Tom Vo¹,
7 Jason Boutte¹, Michael F. Wehner⁹, Angeline G. Pendergrass^{10,11}, Daehyun Kim¹², Zeyu Xue¹³,
8 Andrew T. Wittenberg¹⁴, and John Krasting¹⁴

9

10 ¹ Lawrence Livermore National Laboratory, Livermore, California, USA

11 ² NASA Goddard Space Flight Center, Greenbelt, MD, USA

12 ³ ESSIC, University of Maryland, College Park, MD, USA

13 ⁴ University of California, Davis, Davis, California, USA

14 ⁵ NOAA Pacific Marine Environmental Laboratory, Seattle, Washington, USA

15 ⁶ Monash University, Clayton, Australia

16 ⁷ LOCEAN-IPSL, CNRS-IRD-MNHN-Sorbonne Université, Paris, France

17 ⁸ National Centre for Atmospheric Science-Climate, University of Reading, Reading, UK

18 ⁹ Lawrence Berkeley National Laboratory, Berkeley, California, USA

19 ¹⁰ Department of Earth and Atmospheric Science, Cornell University, Ithaca, New York, USA

20 ¹¹ National Center for Atmospheric Research, Boulder, Colorado, USA

21 ¹² School of Earth and Environmental Sciences, Seoul National University, Seoul, South Korea

22 ¹³ Pacific Northwest National Laboratory, Richland, WA, USA

23 ¹⁴ NOAA Geophysical Fluid Dynamics Laboratory, Princeton, NJ, USA

24 ^a Retired

25

26 Submitted to [Geosci. Model Dev. \(GMD\)](#) in November 2023

27 Revised in December 2023

28

29 *Corresponding to:* Jiwoo Lee (lee1043@llnl.gov)

30 7000 East Ave, Livermore, California 94550, USA

Deleted: [Geoscientific Model Development](#)

32 **Abstract**

33

34 Systematic, routine, and comprehensive evaluation of Earth System Models (ESMs) facilitates benchmarking
35 improvement across model generations and identifying the strengths and weaknesses of different model
36 configurations. By gauging the consistency between models and observations, this endeavor is becoming increasingly
37 necessary to objectively synthesize thousands of simulations contributed to the Coupled Model Intercomparison
38 Project (CMIP) to date. The Program for Climate Model Diagnosis and Intercomparison (PCMDI) Metrics Package
39 (PMP) is an open-source Python software package that provides "quick-look" objective comparisons of ESMs with
40 one another and with observations. The comparisons include metrics of large- to global-scale climatologies, tropical
41 inter-annual and intra-seasonal variability modes such as El Niño-Southern Oscillation (ENSO) and Madden-Julian
42 Oscillation (MJO), extratropical modes of variability, regional monsoons, cloud radiative feedbacks, and high-
43 frequency characteristics of simulated precipitation, including its extremes. The PMP comparison results are produced
44 using all model simulations contributed to CMIP6 and earlier CMIP phases. An important objective of the PMP is to
45 document performance of ESMs participating in the recent phases of CMIP, together with providing version-
46 controlled information for all data sets, software packages, and analysis codes being used in the evaluation process.
47 Among other purposes, this also enables modeling groups to assess performance changes during the ESM development
48 cycle in the context of the error distribution of the multi-model ensemble. Quantitative model evaluation provided by
49 the PMP can assist modelers in their development priorities. In this paper, we provide an overview of the PMP
50 including its latest capabilities, and discuss its future direction.

- Deleted: in the context of
- Deleted: priority
- Deleted: es

54 **1 Introduction**

55 Earth System Models (ESMs) are key tools for projecting climate change and conducting research to enhance
56 our understanding of the Earth system. With the advancements in computing power and the increasing importance of
57 climate projections, there has been an exponential growth of diversity of ESM simulations. During the 1990's, the
58 Atmospheric Model Intercomparison Project (AMIP; Gates, 1992; Gates et al., 1999) was a centralizing activity within
59 the modeling community, which led to the creation of the Coupled Model Intercomparison Project (CMIP; Meehl et
60 al., 1997, 2000, 2007; Covey et al., 2003; Taylor et al., 2012). Since 1989, the Program for Climate Model Diagnosis
61 and Intercomparison (PCMDI) has worked closely with the World Climate Research Programme's (WCRP) Working
62 Group on Coupled Models (WGCM) and Working Group on Numerical Experimentation (WGNE) to design and
63 implement these projects (Potter et al., 2011). The most recent phase of CMIP (CMIP6; Eyring et al., 2016) provides
64 a set of well-defined experiments that most climate modeling centers perform, and subsequently makes results
65 available for a large and diverse community to analyze.

66 Evaluating ESMs is a complex endeavor, given the vast range of climate characteristics across space and
67 time scales. A necessary step involves quantifying the consistency between ESMs with available observations. Climate
68 model performance metrics have been widely used to objectively and quantitatively gauge the agreement between
69 observations and simulations to summarize model behavior with a wide range of climate characteristics. Simple
70 examples include either the model bias or the pattern similarity (correlation) between an observed and simulated field
71 (e.g., Taylor, 2001). With the rapid growth in the number, scale, and complexity of simulations, the metrics have been
72 used more routinely as exemplified by the Intergovernmental Panel on Climate Change (IPCC) Assessment Reports
73 (e.g., Gates et al., 1995; McAvaney et al., 2001; Randall et al., 2007; Flato et al., 2014; Eyring et al., 2021). A few
74 studies have been exclusively devoted to objective model performance assessment using summary statistics. Lambert
75 and Boer (2001) evaluated the first set of CMIP models from CMIP1 using statistics for the large-scale mean climate.
76 Gleckler et al. (2008) identified a variety of factors relevant to model metrics and demonstrated techniques to quantify
77 the relative strengths and weaknesses of the simulated mean climate. Reichler and Kim (2008) attempted to gauge
78 model improvements across the early phases of CMIP. The scope of objective model evaluation has greatly broadened
79 beyond the mean state in recent years (e.g., Gleckler et al., 2016; Eyring et al., 2019), including attempts to establish
80 performance metrics for a wide range of climate variability (e.g., Kim et al., 2009; Sperber et al., 2013; Ahn et al.,
81 2017; Fasullo et al., 2020; Lee et al., 2021b; Planton et al., 2021) and extremes (e.g., Sillmann et al., 2013; Srivastava
82 et al., 2020; Wehner et al., 2020, 2021). Guilyardi et al. (2009) and Reed et al. (2022) emphasized that metrics should
83 be concise, interpretable, informative, and intuitive.

84 With the growth of data size and diversity of ESM simulations, there has been a pressing need for the research
85 community to become more efficient and systematic in evaluating ESMs and documenting their performances. To
86 respond to the need, PCMDI developed the PCMDI Metrics Package (PMP) and released its first version in 2015 (see
87 Code and Data Availability section for all versions). A centralizing goal of the PMP then and now is to quantitatively
88 synthesize results from the archive of CMIP simulations via performance metrics that help characterize the overall
89 agreement between models and observations (Gleckler et al., 2016). For our purposes, "performance metrics" are
90 typically (but not exclusively) well-established statistical measures that quantify the consistency between observed

Deleted: data size and

92 and simulated characteristics. Common examples include a domain average bias, a root-mean-square error (RMSE),
93 a spatial pattern correlation, or others, typically selected depending on the application. Another goal of the PMP is to
94 further diversify the suite of high-level performance tests that help characterize the simulated climate. The results
95 provided by the PMP are frequently used to address two overarching and recurring questions: 1) What are the relative
96 strengths and weaknesses between different models? and 2) How are models improving with further development?
97 Addressing the second question is often referred to as “benchmarking” and this motivates an important emphasis of
98 the effort described in this paper—striving to advance the documentation of all data and results of the PMP in an open
99 and ultimately reproducible manner.

100 In parallel, the current progress towards systematic model evaluation remains dynamic, with evolving
101 approaches and many independent paths being pursued. This has resulted in the development of diversified model
102 evaluation software packages. Examples in addition to the PMP include the ESMValTool (Eyring et al., 2016, 2019,
103 2020; Righi et al., 2020), the Model Diagnostics Task Force (MDTF) Diagnostics package (Maloney et al., 2019;
104 Neelin et al., 2023), the International Land Model Benchmarking (ILAMB) Software System (Collier et al., 2018) that
105 focuses on land surface and carbon cycle metrics, and the International Ocean Model Benchmarking (IOMB) Software
106 System (Fu et al., 2022) that focuses on surface and upper ocean biogeochemical variables. Some tools have been
107 developed with a more targeted focus on a specific subject area, such as the Climate Variability Diagnostics Package
108 (CVDP) that diagnoses climate variability modes (Phillips et al., 2014; Fasullo et al., 2020), and the Analyzing Scales
109 of Precipitation (ASoP) that focuses on analyzing precipitation scales across space and time (Klingaman et al., 2017;
110 Martin et al., 2017; Ordonez et al., 2021). The regional climate community also has actively developed metrics
111 packages such as the Regional Climate Model Evaluation System (RCMES; Lee et al., 2018a; Whitehall et al. 2012).
112 Separately, a few climate modeling centers have developed their own model evaluation packages to assist in their in-
113 house ESM development, e.g., the E3SM Diags (Zhang et al., 2022). There also have been other efforts to enhance
114 the usability of in-situ and field campaign observations in ESM evaluations, such as Atmospheric Radiation
115 Measurement (ARM) GCM Diag (Zhang et al., 2018, 2020) and Earth System Model Aerosol–Cloud Diagnostics
116 (ESMAC Diags; Tang et al., 2022, 2023). While they all have their own scientific priorities and technical approaches,
117 the uniqueness of the PMP is its focus on the objective characterization of the physical climate system as simulated
118 by community models. An important prioritization of the PMP is to advance all aspects of its workflow, in an open,
119 transparent, and reproducible manner, which is critical for benchmarking. The PMP summary statistics characterizing
120 CMIP simulations are version-controlled and made publicly available as a resource to the community.

121 In this paper, we describe the latest update of the PMP and its focus on providing a diverse suite of summary
122 statistics that can be used to construct “quick-look” summaries of ESM performance from simulations made publicly
123 available to the research community, notably CMIP. The rest of the paper is organized as follows. In [Sect. 2](#), we
124 provide a technical description of the PMP and its accompanying reference datasets. In [Sect. 3](#), we describe various
125 sets of simulation metrics that provide an increasingly comprehensive portrayal of physical processes across time
126 scales ranging from hours to centurial. In [Sect. 4](#), we introduce the usage of PMP for model benchmarking. We discuss
127 the future direction and the remaining challenges in [Sect. 5](#) and conclude with a summary in [Sect. 6](#). To assist the
128 reader, the table in Appendix A summarizes the acronyms used in this paper.

- Deleted: section
- Deleted: section
- Deleted: section
- Deleted: section
- Deleted: section

134

135 2 Software package and data description

136 The PMP is a Python-based open-source software framework (https://github.com/PCMDI/pcmdi_metrics)
137 designed to objectively gauge the consistency between ESMs and available observations via well-established statistics
138 such as those discussed in [Sect. 3](#). The PMP has been mainly used for the evaluation of CMIP-participating models.
139 A subset of CMIP experiments, those conducted using the observation forcings such as “Historical” and “AMIP”
140 (Eyring et al., 2016), is particularly well suited for comparing models with observations. The AMIP experiment
141 protocol constrains the simulation with prescribed sea surface temperature (SST), and the “Historical” experiment is
142 conducted using coupled model simulations driven by observed varying natural and anthropogenic forcings. Some of
143 the metrics applicable to these experiments may also be relevant to others (e.g., multi-century coupled control runs
144 called “PiControl” and idealized “4xCO2” simulations that are designed for estimating climate sensitivity).

145 The PMP has been applied to multiple generations of CMIP models in a quasi-operational fashion as new
146 simulations are made available, new analysis methods are incorporated, or new observational data become accessible
147 (e.g., Gleckler et al. 2016; Planton et al., 2021; Lee et al., 2021b; Ahn et al. 2022). Shortly after simulations from the
148 most recent phase of the CMIP (i.e., CMIP6) became accessible, PMP quick-look summaries were provided on the
149 PCMDI’s website (<https://pcmdi.llnl.gov/metrics/>), offering a resource to scientists involved in CMIP or others
150 interested in the evaluation of ESMs. To facilitate this, at PCMDI the PMP is technically linked to the Earth System
151 Grid Federation (ESGF) that is the CMIP data delivery infrastructure (Williams et al., 2016).

152 The primary deliverable of the PMP is a collection of summary statistics. We strive to make the baseline
153 results (raw statistics) publicly available and well-documented, and continue to make advances with this [objective in](#)
154 priority. For our purposes, we are referring to model performance “summary statistics” and “metrics” interchangeably,
155 although in some situations we consider there to be an important distinction. For us, a genuine performance metric
156 constitutes a well-defined and established statistic that has been used in a very specific way (e.g., a particular variable,
157 analysis, and domain) for long-term benchmarking (see [Sect. 4](#)). The distinction between summary statistics and
158 metrics is application-dependent and evolving as the community advances efforts to establish quasi-operational
159 capabilities to gauge ESM performance. Some visualization capabilities described in [Sect. 3](#) are made available
160 through the PMP. Users can also further explore the model data comparisons using their preferred visualization
161 methods or incorporate the results into their own studies from the summary statistics from the PMP. Noting the above,
162 the scope of the PMP is fairly targeted. It is not intended to be “all-purpose”, e.g. by incorporating the vast range of
163 diagnostics used in model evaluation.

164 The PMP is designed to readily work with model output that has been processed using the Climate Model
165 Output Rewriter (CMOR; <https://cmor.llnl.gov/>), which is a software library developed to prepare model output
166 following the CF Metadata Conventions (Hassell et al., 2017; Eaton et al., 2022, <http://cfconventions.org/>) in Network
167 Common Data Form (NetCDF) format. The CMOR is used by most modeling groups contributing to CMIP, ensuring
168 all model output adheres to the CMIP data structures that themselves are based on the CF conventions. It is possible
169 to use the PMP on model output that has not been prepared by CMOR, but this usually requires additional work, e.g.,
170 mapping the data to meet the community standards.

Deleted: Section

Deleted: Section

Deleted: Section

174 For reference datasets, the PMP uses observational products processed to be compliant with the Observations
175 for Model Intercomparison Projects (obs4MIPs; <https://pcmdi.github.io/obs4MIPs/>). The obs4MIPs effort was
176 initiated circa 2010 (Gleckler et al., 2011) to advance the use of the observations in model evaluation and research.
177 Substantial progress has been made in establishing obs4MIPs data standards that technically align with CMIP model
178 output (e.g., Teixeira et al., 2014; Ferraro et al., 2015), with the data products published on the ESGF (Waliser et al.,
179 2020). Obs4MIPs-compliant data were prepared with CMOR, and the data directly available via obs4MIPs are used
180 as PMP reference datasets.

181 The PMP leverages other Python-based open-source tools and libraries such as *xarray* (Hoyer and Hamman,
182 2017), *eofs* (Dawson, 2016), and many others. One of the primary fundamental tools used in the latest PMP version
183 is the Python package, Xarray Climate Data Analysis Tools (xCDAT; Vo et al., 2023; <https://xcdat.readthedocs.io>).
184 The xCDAT is developed to provide a more efficient, robust, and streamlined user experience in climate data analysis
185 when using *xarray* (<https://docs.xarray.dev/>). Portions of the PMP rely on the precursor of the xCDAT, a Python
186 library called Community Data Analysis Tools (CDAT, Williams et al., 2009; Williams, 2014; Doutriaux et al., 2019),
187 which has been fundamental since the early development stages of the PMP. The *xarray* software provides much of
188 the functionality of CDAT (e.g., I/O, indexing, and subsetting). However, it lacks some key climate domain features
189 that have been frequently used by scientists and exploited by the PMP (e.g., regridding, utilization of spatial/temporal
190 bounds for computational operations) which motivated the development of the xCDAT. Completing the transition
191 from CDAT to xCDAT is a technical priority for the next version of PMP.

192 To help advance open and reproducible science, the PMP has been maintained with an open-source policy
193 with accompanying metadata for data reproducibility and reusability. The PMP code is distributed and released with
194 version control. The installation process of PMP is streamlined and user-friendly, leveraging the *Anaconda* distribution
195 and the *conda-forge* channel. By employing *conda* and *conda-forge*, users benefit from a simplified and efficient
196 installation experience, ensuring seamless integration of PMP's functionality with minimal dependencies. This
197 approach not only facilitates a straightforward deployment of the package but also enhances reproducibility and
198 compatibility across different computing environments, thereby facilitating the accessibility and widespread adoption
199 of PMP within the scientific community. The pointer to the installation instructions can be found in the Code and Data
200 Availability section. The PMP's online documentation (http://pcmdi.github.io/pcmdi_metrics/) also includes
201 installation instructions and user demo Jupyter Notebooks. [▲] a database of pre-calculated PMP statistics for all AMIP
202 and Historical simulations in the CMIP archive are also available online. The archive of these statistics, stored as
203 JSON files (Crockford, 2006; Crockford and Morningstar, 2017), includes versioning details for all codes, and
204 dependencies and data that were used for the calculations. These files provide the baseline results of the PMP (see the
205 Code and Data Availability section for details). Advancements in model evaluation along with the number of models
206 and complexity of simulations motivate more systematic documentation of performance summaries. With PMP
207 workflow provenance information being recorded and the model and observational data standards maintained by
208 PCMDI and colleagues, PMP strives to make all its results reproducible.

209

Deleted: We also release a

Deleted: S

212 **3 Current PMP capabilities**

213 The capabilities of the PMP have been expanded beyond its traditional large-scale performance summaries
214 of the mean climate (Gleckler et al., 2008; Taylor, 2001). Various evaluation metrics have been implemented to the
215 PMP for climate variability such as El Niño-Southern Oscillation (ENSO) (Planton et al., 2021; Lee et al., 2021a),
216 extratropical modes of variability (Lee et al., 2019, 2021b), intra-seasonal oscillation (Ahn et al., 2017), monsoons
217 (Sperber and Annamalai, 2014), cloud feedback (Zelinka et al., 2022), and the characteristics of simulated
218 precipitation (Pendergrass et al., 2020; Ahn et al., 2022, 2023) and extremes (Wehner et al., 2020, 2021). These PMP
219 capabilities were built upon model performance tests that have resulted from research by PCMDI scientists and their
220 collaborators. This section will provide an overview of each category of the current PMP evaluation metrics with their
221 usage demonstrations.

222

223 **3.1 Climatology**

224 Mean state metrics quantify how well models simulate observed climatological fields at a large scale, gauged
225 by a suite of well-established statistics such as RMSE, mean absolute error (MAE), and pattern correlation that have
226 been used in climate research for decades. The focus is on the coupled “Historical” and atmospheric-only AMIP (Gates
227 et al., 1999) simulations which are well-suited for comparison with observations. The PMP extracts seasonally and
228 annually averaged fields of multiple variables from large-scale observationally based datasets and results from model
229 simulations. Different obs4MIPs-compliant reference datasets are used depending on the variable examined. When
230 multiple reference datasets are available, one of them is considered as a “default” (e.g., see Table 1) while others are
231 identified as “alternatives”. The default datasets are typically state-of-the-art products, but in general, we lack
232 definitive measures as to which is the most accurate, so the PMP metrics are routinely calculated with multiple
233 products so that it can be determined what difference the selection of alternative observations makes to judgment made
234 about model fidelity. The suite of mean climate metrics (all area weighted) includes spatial and spatiotemporal RMSE,
235 centered spatial RMSE, spatial-mean bias, spatial standard deviation, spatial pattern correlation, and spatial and
236 spatiotemporal MAE of the annual or seasonal climatological time-mean (Gleckler et al., 2008). Often, a space-time
237 statistic is used that gauges both the consistency of the observed and simulated climatological pattern as well as its
238 seasonal evolution (see Eq. 1 from Gleckler et al., 2008). By default, results are available for selected large-scale
239 domains, including: “Global”, “Northern Hemisphere (NH) Extratropics” (30°N-90°N), “Tropics” (30°S-30°N), and
240 “Southern Hemisphere (SH) Extratropics” (30°S-90°S). For each domain, results can also be computed for the land
241 and ocean, land only, or ocean only. These commonly used domains highlight the application of the PMP mean climate
242 statistics at large to global scales, but we note that PMP allows users to define their own domains of interest, including
243 at regional scales. Detailed instructions can be found on the PMP’s online documentation
244 (http://pcmdi.github.io/pcmdi_metrics).

245 Although the primary deliverable of the PMP is the metrics, the PMP results can be visualized in various
246 ways. For individual fields, we often first plot Taylor Diagrams, a polar plot leveraging the relationship between the
247 centered RMSE, the pattern correlation, and the observed and simulated standard deviation (Taylor, 2001). The Taylor
248 Diagram has become a standard plot in the model evaluation workflow across modeling centers and research

Deleted: se

250 communities (see Sect. 5). To interpret results across CMIP models for many variables, we routinely construct
251 normalized Portrait Plots or Gleckler Plots (Gleckler et al., 2008) that provide a quick-look examination of the
252 strengths and weaknesses of different models. For example, in Figure 1, the PMP results display quantitative
253 information of simulated seasonal climatologies of various meteorological model variables via a normalized global
254 spatial RMSE (Gleckler et al., 2008). Variants of this plot have been widely used for presenting model evaluation
255 results, for example, in the IPCC Fifth (Flato et al., 2014, Figures 9.7, 9.12, and 9.37) and Sixth Assessment Reports
256 (Eyring et al., 2021, Chapter 3, Figure 3.42). Because the error distribution across models is variable dependent, the
257 statistics are often normalized to help reveal differences, in this case via the median RMSE across all models (see
258 Gleckler et al. 2008 for more details). This normalization enables a common color scale to be used for all statistics on
259 the Portrait Plot, highlighting the relative strengths and weaknesses of different models. In this example (Fig. 1), an
260 error of -0.5 indicates that a model's error is 50% smaller than the typical (median) error across all models, whereas
261 an error of 0.5 is 50% larger than the typical error in the multi-model ensemble. In many cases, the horizontal bands
262 in the Gleckler plots show that simulations from a given modeling center have similar error structures relative to the
263 multi-model ensemble.

264 The Parallel Coordinate Plot (Inselberg, 1997, 2008, 2016; Johansson and Forsell, 2016) that retains the
265 absolute value of the error statistics is used to complement the Portrait plot. Some previous studies have utilized
266 Parallel Coordinate Plots for analyzing climate model simulations (e.g., Steed et al., 2012; Wong et al., 2014; Wang
267 et al., 2017), but to date, only a few studies have applied it to collective multi-ESM evaluations (e.g., see Fig. 7 of
268 Boucher et al., 2020). In the PMP, we generally construct Parallel Coordinate Plots using the same data as in a portrait
269 plot. However, a fundamental difference is that metrics values can be more easily scaled to highlight absolute values
270 rather than the normalized relative results of the portrait plot. In this way, the Portrait and Parallel Coordinate plots
271 complement one another, and in some applications, it can be instructive to display both. Figure 2 shows the
272 spatiotemporal RMSE, defined as the temporal average of spatial RMSE calculated in each month of the annual cycle,
273 of CMIP5 and CMIP6 models in the format of Parallel Coordinate Plot. Each vertical axis represents a different scalar
274 measure gauging a distinct aspect of model fidelity. While polylines are frequently used to connect data points from
275 the same source (i.e., metric values from the same model, in our case) in Parallel Coordinate Plots, we display results
276 from each model using an identification symbol to reduce visual clutter on the plot and help identify outlier models.
277 In the example of Fig. 2, each vertical axis is aligned with the median value midway through its max/min range scale.
278 Thus, for each axis, the models in the lower half of the plot perform better than the CMIP5-CMIP6 multi-model
279 median, while in the upper half, the opposite is true. For each vertical axis that is for a different model variable, we
280 have added violin plots (Hintze and Nelson, 1998) to show probability density functions representing the distributions
281 of model performance obtained from CMIP5 (shaded in blue, left side of the axis) and CMIP6 (shaded in orange, right
282 side of the axis). Medians of each CMIP5 and CMIP6 group are highlighted using polylines, which indicates that the
283 RMSE is reduced in CMIP6 relative to CMIP5 in general for the majority of the subset of model variables.

284

Deleted: Section

286 **3.2 El Niño-Southern Oscillation**

287 The El Niño-Southern Oscillation (ENSO) is Earth's dominant interannual mode of climate variability, which
288 impacts global climate via both regional oceanic effects and far-reaching atmospheric teleconnections (McPhaden et
289 al., 2006, 2020). In response to increasing interest in a community approach to ENSO evaluation in models (Bellenger
290 et al., 2014), the International Climate and Ocean Variability, Predictability and Change (CLIVAR) Research Focus
291 on ENSO in a Changing Climate, together with the CLIVAR Pacific Region Panel, developed the CLIVAR ENSO
292 Metrics Package (Planton et al., 2021) which is now utilized within the PMP. The ENSO metrics used to
293 assess/evaluate the models are grouped into three categories: *Performance* (i.e., background climatology and basic
294 ENSO characteristics), *Teleconnections* (ENSO's worldwide teleconnections), and *Processes* (ENSO's internal
295 processes and feedback). Planton et al. (2021) found that CMIP6 models generally outperform CMIP5 models in
296 several ENSO metrics in particular for those related to tropical Pacific seasonal cycles and ENSO teleconnections.
297 This effort is discussed in more detail in Planton et al. (2021), and detailed descriptions of each metric in the package
298 are available in the ENSO Package online open-source code repository on its GitHub Wiki pages (see
299 https://github.com/CLIVAR-PRP/ENSO_metrics/wiki).

300 Figure 3 demonstrates the application of the ENSO metrics to CMIP6, showing the magnitudes of inter-
301 model and inter-ensemble spreads, along with observational uncertainty varying across metrics. For a majority of the
302 ENSO Performance metrics model error and inter-model spread are substantially larger than observational uncertainty
303 (Figs. 3a-n). This highlights the systematic biases like the double intertropical convergence zone (ITCZ) (Fig. 3a) that
304 are persisting through CMIP phases (Tian and Dong, 2020). Similarly, ENSO Processes metrics (Figs. 3t-w) indicate
305 large errors in the feedback loops generating SST anomalies, indicating a different balance of processes in the model
306 and in the reference and possibly compensating errors (Bayr et al., 2019, Guilyardi et al. 2020). In contrast, for ENSO
307 Teleconnection metrics, the observational uncertainty is substantially larger, thus challenging validation of model
308 error (Figs. 3o-r). For some metrics, such as the ENSO duration (Fig. 3f), the ENSO Asymmetry metric (Fig. 3i), and
309 the Ocean driven SST metric (Fig. 3s), there are larger inter-ensemble spreads than the inter-model spreads. From
310 such results, Lee et al. (2021a) examined the inter-model and inter-member spread of these metrics from the large
311 ensembles available from CMIP6 and the US CLIVAR Large Ensemble Working Group. They argued that to robustly
312 characterize baseline ENSO characteristics and physical processes, larger ensemble sizes are needed, compared to
313 existing state-of-the-art ensemble projects. By applying the ENSO metrics to historical and piControl simulations of
314 CMIP6 via the PMP, Planton et al. (2023) developed equations based on statistical theory to estimate the required
315 ensemble size for a user-defined uncertainty range.

316
317 **3.3 Extratropical Modes of Variability**

318 The PMP includes objective measures of the pattern and amplitude of extratropical modes of variability from
319 PCMDI's research, which has expanded beyond its traditional large-scale performance summaries to include
320 interannual variability, considering increasing interest in setting an objective approach for the collective evaluation of
321 multiple modes. Extratropical modes of variability (ETMoV) metrics in the PMP were developed by Lee et al. (2019a)
322 that stem from earlier works (e.g., Stoner et al., 2009; Phillips et al., 2014). Lee et al. (2019a) illustrated a challenge

323 when evaluating modes of variability using the traditional empirical orthogonal functions (EOF). In particular, when
324 a higher-order EOF of a model more closely corresponds to a lower-order observationally based EOF (or vice versa),
325 it can significantly affect conclusions drawn about model performance. To circumvent this issue in evaluating the
326 interannual variability modes, Lee et al. (2019a) used the Common Basis Function (CBF) approach that projects the
327 observed EOF pattern onto model anomalies. This approach has been previously applied for the evaluation of
328 intraseasonal variability modes (Sperber, 2004; Sperber et al., 2005). In the PMP, the CBF approach is taken as a
329 default method, and the traditional EOF approach is also enabled as an option for the ETMoV metrics calculations.

330 The ETMoV metrics in the PMP measure simulated patterns and amplitudes of ETMoV, and quantify their
331 agreement with observations (e.g., Lee et al., 2019a, 2021b). The PMP's ETMoV metrics evaluate 5 atmospheric
332 modes – the Northern Annular Mode (NAM), North Atlantic Oscillation (NAO), Pacific North America pattern
333 (PNA), North Pacific Oscillation (NPO), and Southern Annular Mode (SAM), and 3 ocean modes diagnosed by the
334 variance of sea-surface temperature – Pacific Decadal Oscillation (PDO), North Pacific Gyre Oscillation (NPGO),
335 and Atlantic Multi-decadal Oscillation (AMO). The AMO is included for experimental purposes, considering the
336 significant uncertainty in detecting the AMO (Deser and Philips 2021; Zhao et al., 2022). The amplitude metric,
337 defined as the ratio of standard deviations of the model and observed principal components, has been used to examine
338 the evolution of the performance of models across different CMIP generations (Fig. 4). Green shading predominates,
339 indicating where the simulated amplitude of variability is similar to observations. In some cases, such as for SAM in
340 September-October-November (SON), the models overestimate the observed amplitude.

341 The PMP's ETMoV metrics have been used in several model evaluation studies. For example, Orbe et al.
342 (2020) analyzed models from U.S. climate modeling groups including the U.S. Department of Energy (DOE), National
343 Aeronautics and Space Administration (NASA), National Center for Atmospheric Research (NCAR), and National
344 Oceanic and Atmospheric Administration (NOAA), where they found that the improvement in the ETMoV
345 performance is highly dependent on mode and season, when comparing across different generations of those models.
346 Sung et al. (2021) examined the performance of models run at the Korea Meteorological Administration (K-ACE and
347 UKESM1) in reproducing ETMoVs from their Historical simulations, and concluded that these models reasonably
348 capture most ETMoVs. Lee et al. (2021b) collectively evaluated ~130 models from CMIP3, 5, and 6 archive databases
349 using their ~850 Historical and ~300 AMIP simulations, where they found the spatial pattern skill improved in CMIP6
350 compared to CMIP5 or CMIP3 for most modes and seasons, while the improvement in amplitude skill is not clear.
351 Arcodia et al. (2023) used the PMP to derive PDO and AMO to investigate their role in decadal variability of
352 subseasonal predictability of precipitation over the western coast of North America and concluded that no significant
353 relationship was found.

354 355 **3.4 Intraseasonal Oscillation**

356 The PMP has implemented metrics for the Madden-Julian Oscillation (MJO; Madden and Julian, 1971, 1972,
357 1994). The MJO is the dominant mode of tropical intraseasonal variability, characterized by a pronounced eastward
358 propagation of large-scale atmospheric circulation coupled with convection with a typical periodicity of 30-60 days.

Deleted: _

360 Selected metrics from the MJO diagnostics package, developed by the CLIVAR MJO Working Group (Waliser et al.,
361 2009), have been implemented in the PMP following Ahn et al. (2017).

362 We have particularly focused on metrics for the MJO propagation: East/West power Ratio (EWR) and East
363 power normalized by Observation (EOR). The EWR is proposed by Zhang and Hendon (1997) which is defined as
364 the ratio of the total spectral power over the MJO band (eastward propagating, wavenumber 1-3 and period of 30-60
365 days) to that of its westward propagating counterpart in the wavenumber-frequency power spectra. The EWR metric
366 has been widely used in the community, to examine the robustness of the eastward propagating feature of the MJO
367 (e.g., Hendon et al., 1999; Lin et al., 2006; Kim et al., 2009; Ahn et al., 2017). The EOR is formulated by normalizing
368 a model's spectral power within the MJO band by the corresponding observed value. Ahn et al. (2017) showed EWRs
369 and EORs of the CMIP5 models. Using daily precipitation, the PMP calculates EWR and EOR separately for boreal
370 winter (November to April) and boreal summer (March to October). We apply the frequency-wavenumber
371 decomposition method to precipitation from observations (GPCP-based; 1997-2010) and the CMIP5 and CMIP6
372 Historical simulations for 1985-2004. For disturbances with wavenumbers 1-3 and frequencies corresponding to 30-
373 60 days, it is clear in observations that the eastward propagating signal dominates over its westward propagating
374 counterpart with an EWR value of about 2.49 (Fig. 5a). Figure 5b shows the wavenumber-frequency power spectrum
375 from CMIP5 IPSL-CM5B-LR as an example, which has an EWR value that is comparable to the observed value.

376 Figure 6 shows the EWR from individual models' multiple ensemble members and their average. The average
377 EWR of the CMIP6 model simulations is more realistic than that of the CMIP5 models. Interestingly, a substantial
378 spread exists across models and also among ensemble members of a single model. For example, while the average
379 EWR value for the CESM2 ensemble is 2.47 (close to 2.49 from the GPCP observations), the EWR values of the
380 individual ensemble members range from 1.87 to 3.23. Kang et al. (2020) suggested that the ensemble spread in the
381 propagation characteristics of the MJO can be attributed to the differences in the moisture mean state, especially its
382 meridional moisture gradient. A cautionary note should be given to the fact that the MJO frequency and wavenumber
383 windows are chosen to capture the spectral peak in observations. Thus, while the EWR provides an initial evaluation
384 of the propagation characteristics of the observed and simulated MJO, it is instructive to look at the frequency-
385 wavenumber spectra, as in some cases the dominant periodicity and wavenumber in a model may be different than in
386 observations. It is worthwhile to note that the PMP can be used to obtain EWR and EOR of other daily variables for
387 MJO analysis, such as outgoing longwave radiation (OLR) or zonal wind at 850 hPa (U-850) or 250 hPa (U-250), as
388 shown in Ahn et al. (2017).

389

390 3.5 Monsoons

391 Based on the work of Sperber and Annamalai (2014), skill metrics in the PMP quantify how well models
392 represent the onset, decay, and duration of regional monsoons. From observations and Historical simulations, the
393 climatological pentad ~~data~~ of precipitation are area-averaged for six monsoon domains: All-India Rainfall, Sahel, Gulf
394 of Guinea, North American Monsoon, South American Monsoon, and Northern Australia (Fig. 7). For the domains in
395 the Northern Hemisphere, the 73 climatological pentads run from January to December, while for the domains in the
396 Southern Hemisphere, the pentads run from July to June. For each domain, the precipitation is accumulated at each

Deleted: s

Deleted: -related

Deleted: , as seen in

400 subsequent pentad and then divided by the total precipitation to give the fractional accumulation of precipitation as a
401 function of pentad. Thus, the annual cycle behavior is evaluated irrespective of whether a model has a dry or wet bias.
402 Except for the Gulf of Guinea, the onset and decay of monsoon occur for a fractional accumulation of 0.2 and 0.8,
403 respectively. Between these fractional accumulations, the accumulation of precipitation is nearly linear as the monsoon
404 season progresses. Comparison of the simulated and observed onset, duration, and decay are presented in terms of the
405 difference in the pentad index obtained from the model and observations (i.e., model minus observations). Therefore,
406 negative values indicate that the onset or decay in the model occurs earlier than in observations, while positive values
407 indicate the opposite. For duration, negative values indicate that for the model it takes fewer pentads to progress from
408 onset to decay compared to observations (i.e., the simulated monsoon period is too short), while positive values
409 indicate the opposite.

410 For CMIP5, we find systematic errors in the phase of the annual cycle of rainfall. The models are delayed in
411 the onset of summer rainfall over India, the Gulf of Guinea, and the South American Monsoon, with early onset
412 prevalent for the Sahel and the North American Monsoon. The lack of consistency in the phase error across all domains
413 suggests that a “global” approach to the study of monsoons may not be sufficient to rectify the regional differences.
414 Rather, regional process studies are necessary for diagnosing the underlying causes of the regionally specific
415 systematic model biases over the different monsoon domains. Assessment of the monsoon fidelity in CMIP6 models
416 using the PMP is in progress.

417

418 *3.6 Cloud feedback and mean-state*

419 Uncertainties in cloud feedback are the primary driver of model-to-model differences in climate sensitivity
420 – the global temperature response to a doubling of atmospheric CO₂. Recently, an expert synthesis of several lines of
421 evidence spanning theory, high-resolution models, and observations was conducted to establish quantitative
422 benchmark values (and uncertainty ranges) for several key cloud feedback mechanisms. The assessed feedbacks are
423 those due to changes in high-cloud altitude, tropical marine low-cloud amount, tropical anvil cloud area, land cloud
424 amount, middle latitude marine low-cloud amount, and high latitude low-cloud optical depth. The sum of these six
425 components yields the total assessed cloud feedback, which is part of the overall radiative feedback that fed into the
426 Bayesian calculation of climate sensitivity in Sherwood et al. (2020). Zelinka et al. (2022) estimated these same
427 feedback components in climate models and evaluated them against the expert-judgment values determined in
428 Sherwood et al. (2020), ultimately deriving a root mean square error metric that quantifies the overall match between
429 each model’s cloud feedback and those determined through expert judgment.

430 Figure 8 shows the model-simulated values for each individual feedback computed in *amip-p4K* simulations
431 as part of CMIP5 and CMIP6 alongside the expert judgment values. Each model is color-coded by its equilibrium
432 climate sensitivity (determined using *abrupt-4xCO2* simulations as described in Zelinka et al., 2020), and the values
433 from an illustrative model (GFDL-CM4) are highlighted. Among the key results apparent from this figure is that
434 models typically underestimate the strength of both positive tropical marine low-cloud feedback and the negative anvil
435 cloud feedback relative to the central expert assessed value. The sum of all six assessed feedback components is

436 positive in all but two models, with a multimodel mean value that is close to the expert-assessed value, but exhibits
437 substantial intermodel spread.

438 In addition to evaluating the ability of models to match the assessed cloud feedback components, Zelinka et
439 al. (2022) investigated whether models with less erroneous mean-state clouds tend to have smaller errors in their
440 overall cloud feedback RMSE. This involved computing the mean-state cloud property error metric developed by
441 Klein et al. (2013). This error metric quantifies the spatiotemporal error in climatological cloud properties for clouds
442 with optical depths greater than 3.6, weighted by their net top-of-atmosphere (TOA) radiative impact. The
443 observational baseline against which the models are compared comes from the International Satellite Cloud
444 Climatology Project H-series Gridded Global (ISCCP HGG) dataset (Young et al., 2018). Zelinka et al. (2022) showed
445 that models with smaller mean-state cloud errors tend to have stronger but not necessarily better (less erroneous) cloud
446 feedback, which suggests that improving mean-state cloud properties does not guarantee improvement in the cloud
447 response to warming. However, the models with the smallest errors in cloud feedback tend also to have less erroneous
448 mean-state cloud properties, and no models with poor mean-state cloud properties have feedback in good agreement
449 with expert judgment.

450 The PMP implementation of this code computes cloud feedback by differencing fields from *amip-p4K* and
451 *amip* experiments and normalizing by the corresponding global mean surface temperature change rather than from
452 differencing *abrupt-4xCO2* and *piControl* experiments and computing feedback via regression (as was done in Zelinka
453 et al., 2022). This choice is made to reduce the computational burden and also because cloud feedbacks derived from
454 these simpler atmosphere-only simulations have been shown to closely match those derived from fully coupled
455 quadrupled CO₂ simulations (Qin et al., 2022). The code produces figures in which the user-specified model results
456 are highlighted and placed in the context of the CMIP5 and CMIP6 multi-model results (e.g., Fig. 8).

Deleted: Figure

458 3.7 Precipitation

459 Recognizing the importance of accurately simulating precipitation in ESMs and a lack of objective and
460 systematic benchmarking for it, and motivated by discussions with WGNE and WGCM working groups of WCRP,
461 the DOE has initiated an effort to establish a pathway to help modelers gauge improvement (U.S. DOE, 2020). The
462 2019 DOE workshop “Benchmarking Simulated Precipitation in Earth System Models” generated two sets of
463 precipitation metrics: *baseline* and *exploratory* metrics (Pendergrass et al., 2020). In the PMP, we have focused on
464 implementing the *baseline* metrics for benchmarking simulated precipitation. In parallel, a set of *exploratory* metrics
465 that could be added to metrics suites including PMP in the future was illustrated by Leung et al. (2022) to extend the
466 evaluation scope to include process-oriented and phenomena-based diagnostics and metrics.

467 The *baseline* metrics gauge the consistency between ESMs and observations, focusing on the holistic set of
468 observed rainfall characteristics (Fig. 9). For example, the spatial distribution of mean state precipitation and seasonal
469 cycle are outcomes of the PMP’s Climatology metrics (described in Sect. 3.1), which provides collective evaluation
470 statistics such as RMSE, standard deviation, and pattern correlation over various domains (e.g., global, NH and SH
471 extratropics, and tropics, with each domain as a whole, and over land and ocean, in separate). Evaluation of
472 precipitation variability across many timescales with PMP is documented in Ahn et al. (2022); we summarize some

Deleted: Section

Deleted: Tropics

476 of the findings here. The precipitation variability metric measures forced (diurnal and annual cycles) and internal
477 variability across timescales (subdaily, synoptic, subseasonal, seasonal, and interannual) in a framework based on
478 power spectra of 3-hourly total and anomaly precipitation. Overall, CMIP5 and CMIP6 models underestimate the
479 internal variability, which is more pronounced in the higher frequency variability, while they overestimate the forced
480 variability (Fig. 10). For the diurnal cycle, PMP includes metrics from Covey et al. (2016). Additionally, the intensity
481 and distribution of precipitation are assessed following Ahn et al. (2023). Extreme daily precipitation indices and their
482 20-year return values are calculated using a non-stationary Generalized Extreme Value statistical method. From the
483 CMIP5 and CMIP6 historical simulations we evaluate model performance of these indices and their return values in
484 comparison with gridded land-based daily observations. Using this approach, Wehner et al. (2020) found that at
485 models' standard resolutions, no meaningful differences were found between the two generations of CMIP models.
486 Wehner et al. (2021) extended the evaluation of simulated extreme precipitation to seasonal 3-hourly precipitation
487 extremes produced by available HighResMIP models and concluded that the improvement is minimal with the models'
488 increased spatial resolutions. They also noted that the order of operations of regridding and calculating extremes
489 affects the ability of models to reproduce observations. Drought metrics developed by Xue and Ullrich (2021) are not
490 implemented in PMP directly, but are wrapped by the Coordinated Model Evaluation Capabilities (CMEC; Ordonez
491 et al. 2021), which is a parallel framework for supporting community-developed evaluation packages. Together, these
492 metrics provide a streamlined workflow for running the entire baseline metrics via the PMP and CMEC that is ready
493 for use by operational centers and in the CMIP7.

494

495 *3.8 Relating metrics to underlying diagnostics*

496 Considering the extensive collection of information generated from the PMP, efforts have supported
497 improved visualizations of metrics using interactive graphic user interfaces. These capabilities can facilitate the
498 interpretation and synthesis of vast amounts of information associated with the diverse metrics and the underlying
499 diagnostics from which they were derived. Via the interactive navigation interface, we can explore the underlying
500 diagnostics behind the PMP's summary plots. On the PCMDI website, we provide interactive graphical interfaces to
501 enable navigating the supporting plots to the underlying diagnostics of each model's ensemble members and their
502 average. For example, on the interactive mean climate plots (https://pcmdi.llnl.gov/metrics/mean_clim/), hovering the
503 mouse cursor over a square or triangle in the Portrait Plot, or over the markers or lines in the Parallel Coordinate Plot,
504 reveals the diagnostic plot from which the metrics were generated. It allows the user to toggle between several metrics
505 (e.g., RMSE, bias, and correlation) and regions (e.g., global, Northern/Southern Hemisphere, and Tropics), along with
506 relevant provenance information. Users can click on the interactive plots to get dive-down diagnostics information for
507 the model of interest which provides detailed analysis to better understand how the metric was calculated. As with the
508 PMP's mean climate metrics output, we currently provide interactive summary graphics for ENSO
509 (<https://pcmdi.llnl.gov/metrics/enso/>), extratropical modes of variability
510 (https://pcmdi.llnl.gov/metrics/variability_modes/), monsoon (<https://pcmdi.llnl.gov/metrics/monsoon/>), MJO
511 (<https://pcmdi.llnl.gov/metrics/mjo/>), and precipitation benchmarking (<https://pcmdi.llnl.gov/metrics/precip/>). We
512 plan to expand this capability to other metrics in the PMP, such as the cloud feedback analysis. The majority of the

Deleted: s

514 PMP's interactive plots have been developed using Bokeh (<https://bokeh.org/>), a Python data visualization library that
515 enables the creation of interactive plots and applications for web browsers.

516

517 **4 Model Benchmarking**

518 While the PMP originally focused on evaluating multiple models (e.g., Gleckler et al., 2008), in parallel there
519 has been increasing interest from model developers and modeling centers to leverage the PMP to track performance
520 evolution in the model development cycle, as discussed in Gleckler et al. (2016). For example, metrics from the PMP
521 have been used to document performance of ESMs developed in the U.S. DOE Exascale Earth System Model (E3SM;
522 Caldwell et al., 2019; Golaz et al., 2019; Rasch et al., 2019; Hannah et al., 2021; Tang et al., 2021), NOAA
523 Geophysical Fluid Dynamics Laboratory (GFDL; Zhao et al., 2018), Institut Pierre-Simon Laplace (IPSL; Boucher et
524 al., 2020; Planton et al., 2021), National Institute of Meteorological Sciences-Korea Meteorological Administration
525 (NIMS-KMA; Sung et al., 2021), University of California, Los Angeles (Lee et al., 2019b), and the Community
526 Integrated Earth System Model (CIesm) project (Lin et al., 2020).

527 To make the PMP more accessible and useful for modeling groups, efforts are underway to broaden workflow
528 options. Currently, a typical application involves computing a particular class of performance metrics (e.g., mean
529 climate) for all CMIP simulations available via ESGF. To facilitate the ability of modeling groups to routinely use the
530 PMP during their development process, we are working to provide a customized workflow option to run all the PMP
531 metrics more seamlessly on a single model, and to compare these results with a database of PMP results obtained from
532 CMIP simulations (see Code and Data Availability section). Via the PMP-documented and pre-calculated metrics
533 from simulations in the CMIP archive, it is possible to readily incorporate CMIP results into the assessment of new
534 simulations, without retrieving all CMIP simulations and recomputing the results. The resulting quick-look feedback
535 can highlight model improvement (or deterioration) and can assist in determining development priorities or in the
536 selection of a new model version.

537 As an example, here, we show PMP results obtained from GFDL-CM3 from CMIP5 and GFDL-CM4 from
538 CMIP6, for a demonstration of using the Taylor Diagram to compare versions of a given model (Fig. 11). One
539 advantage of the Taylor Diagram is that it collectively represents three statistics (i.e., centered RMSE, standard
540 deviation, and correlation) in a single plot (Taylor, 2001), which synthesizes the performance intercomparison of
541 multiple models (or different versions of a model). In this example, four variables were selected to summarize
542 performance evolution (shown by arrows) in multiple seasons. Except for boreal winter, both model versions are
543 nearly identical in terms of net TOA radiation, however in all seasons the longwave cloud radiative effect is clearly
544 improved in the newer model version. The TOA flux improvements likely contributed to the precipitation
545 improvements, by improving the balances of radiative cooling and latent heating. The improvement in the newer
546 model version is consistent with that documented by Held et al., (2019) and evident via the arrow directions pointing
547 to the observational reference point.

548 Parallel Coordinate Plots can also be used to summarize the comparison of two simulations for their
549 performance. In Fig 12, we demonstrate the comparison of selected metrics: the mean climate (see [Sect. 3.1](#)), ENSO
550 ([Sect. 3.2](#)), and ETMoV ([Sect. 3.3](#)). To facilitate comparison of a subset of models, a few models can be selected and

Deleted: Section
Deleted: Section
Deleted: Section

554 highlighted as connected lines across individual vertical axes on the plot. A proposed application of it from PMP is to
555 select two models or two versions of a model to contrast their performance (solid lines) against the backdrop of results
556 from other models, shown as violin plots for the distribution of statistics from other models on each vertical axis. In
557 this example, we contrast the performance of two GFDL models: GFDL-CM3 and GFDL-CM4. Fig 12a is a modified
558 version of Figure 2 that is designed to highlight the difference in performance more efficiently. Each vertical axis
559 indicates performance for each metric defined for climatology of variables (i.e., temporally averaged spatial RMSE
560 of annual cycle climatology patterns, Fig. 12a), ENSO characteristics (Fig. 12b), or interannual variability mode
561 obtained from seasonal or monthly averaged time series (Fig. 12c). It is shown that GFDL-CM4 is superior to GFDL-
562 CM3 for most cases across selected metrics (downward arrows in green) while inferior for a few cases (upward arrows
563 in red), which is consistent with previous findings (Held et al., 2019; Planton et al., 2021; Chen et al., 2021). Such
564 applications of the Parallel Coordinate Plot can enable quick overall assessment and tracking of the ESM performance
565 evolution during its development cycle. More examples showing other models are available in the Supplementary
566 material (Figs. S1 to S3).

567 It is worth noting that there have been efforts to coalesce objective model evaluation concepts used in the
568 research community (e.g., Knutti et al., 2010). However, the field continues to evolve rapidly with definitions still
569 being debated and finessed. Via the PMP, we produce hundreds of summary statistics, enabling a broad net to be cast
570 in the objective characterization of a simulation, at times helping modelers identify previously unknown deficiencies.
571 For benchmarking, efforts are underway to establish a more targeted path which likely involves a consolidated set of
572 carefully selected metrics.

574 **5 Discussion**

575 Efforts are underway to include new metrics into the PMP to advance the systematic objective evaluation of
576 ESMs. For example, in coordination with the World Meteorological Organization (WMO)'s WGNE MJO Task Force,
577 additional candidate MJO metrics for PMP inclusion have been identified to facilitate more comprehensive
578 assessments of the MJO. Implementation of metrics for MJO amplitude, periodicity, and structure into the PMP is
579 planned. An ongoing collaboration with NCAR aims to incorporate metrics related to the upper atmosphere,
580 specifically the Quasi-Biennial Oscillation (QBO) and QBO-MJO metrics (e.g. Kim et al., 2020). We also have plans
581 to grow the scope of PMP beyond its traditional atmospheric realm, for example including the ocean and polar regions
582 through collaboration with the U.S. DOE's project entitled High Latitude Application and Testing of ESMs (HiLAT,
583 <https://www.hilat.org/>). In addition, the PMP framework is also well poised to contribute to high-resolution climate
584 modeling activities, such as the High-Resolution Model Intercomparison Project (HighResMIP; Haarsma et al., 2016)
585 and the DYNamics of the Atmospheric general circulation Modeled On Non-hydrostatic Domains (DYAMOND;
586 Stevens et al., 2019). This motivates the development of specialized metrics for high-resolution models, targeting the
587 simulation features enabled by high-resolution models. Another potential avenue for the PMP involves leveraging
588 Machine Learning (ML) techniques, and other state-of-the-art data science techniques being used for process-oriented
589 ESM evaluation works (e.g., Nowack et al., 2020; Labe and Barnes, 2022; Dalelane et al., 2023). Applications of ML
590 detection, such as for storms using TempestExtremes (Ullrich and Zarzycki 2017; Ullrich et al., 2021) and fronts (e.g.,

591 Biard and Kunkel, 2019), can enable additional specialized storm metrics for high-resolution simulations. For
592 convection-permitting models, yet more storm metrics can be applied such as Mesoscale convective systems.
593 Atmospheric blocking metrics and atmospheric river evaluation metrics using the ML pattern detection capabilities in
594 the latest TempestExtremes (Ullrich et al., 2021) are currently under development to be implemented into the PMP.
595 These example enhancements of the PMP are indicative of an increasing priority to target regional simulation
596 characteristics. With a deliberate emphasis on processes intrinsic to specific regions, this may lead to enabling
597 potential applications of the PMP within the regional climate modeling activities such as the Coordinated Regional
598 Downscaling Experiment (CORDEX; Gutowski Jr. et al., 2016).

599 The comprehensive database of PMP results offers a resource for exploring the range of structural errors in
600 CMIP class models and their interrelationships. For example, examination of cross-metric relationships between
601 mean-state and variability biases can shed additional light on the propagation of errors (e.g., Kang et al., 2020; Lee et
602 al., 2021b). There continues to be interest in ranking models for specific applications (e.g., Ashfaq et al., 2022;
603 Goldenson et al., 2023; Longmate et al., 2023; Papalexiou et al., 2020; Singh and AchutaRao, 2020) or to “move
604 beyond one model one vote” in multi-model analysis to reduce uncertainties in the spread of multi-model projections
605 (e.g., Knutti, 2010; Knutti et al., 2017; Sanderson et al., 2017; Herger et al., 2018; Hausfather et al., 2022; Merrifield
606 et al., 2023). While we acknowledge potential interests in using the results of the PMP or equivalent to rank models
607 or identify performance outliers (e.g., Sanderson and Wehner, 2017), we believe the many challenges associated with
608 model weighting are application dependent, and thus leave it up to users of the PMP to make those judgments.

609 In addition to the scientific challenges associated with diversifying objective summaries of model
610 performance, there is potential to leverage rapidly evolving technologies, including new open-source tools and
611 methods available to scientists. We expect that the ongoing PMP code modernization effort to fully adapt the xCDAT
612 and xarray will facilitate greater community involvement. As the PMP evolves with these technologies we will
613 continue to maintain rigor in the calculation of statistics for the PMP metrics, for example by incorporating the latest
614 advancements in the field. A prominent example in the objective comparison of models and observations involves the
615 methodology of horizontal interpolation, and in future versions of the PMP we are planning a more stringent
616 conservation method (Taylor, 2024). To improve the clarity of key messages from multivariate PMP metrics data, we
617 will consider implementing the advances in high-dimensional data visualization, e.g., the circular plot discussed in
618 Lee et al. (2018b) and variations of Parallel Coordinate Plots proposed in this paper and by Hassan et al. (2019) and
619 Lu et al. (2020).

620 Current progress towards systematic model evaluation is exemplified by the diversity of tools being
621 developed (e.g., the PMP, ESMValTool, MDTF, ILAMB, IOMB, and other packages). Each of these tools has its own
622 scientific priorities and technical approaches. We believe that this diversity has made, and will continue to make,
623 the model evaluation process even more comprehensive and successful. The fact that there is some overlap in a few cases
624 is advantageous because it enables the cross-verification of results, which is particularly useful in more complex
625 analyses. Despite possible advantages, having no single best or widely accepted approach for the community to follow,
626 does introduce complexity to the coordination of model evaluation. To facilitate the collective usage of individual
627 evaluation tools, the CMEC has initiated the development of a unified code base that technically coordinates the

Deleted: 3

629 operation of distinct but complementary tools (Ordonez et al. 2021). Currently, the PMP, ILAMB, MDTF, and ASoP
630 have become CMEC-compliant by adopting common interface standards that define how evaluation tools interact
631 with observational data and climate model output. We expect that CMEC can also help the model evaluation
632 community to establish standards for archiving the metrics output, much as the community did for the conventions to
633 describe climate model data (e.g., CMIP application of CF Metadata Conventions (<http://cfconventions.org/>); Hassell
634 et al., 2017; Eaton et al., 2022).

635 **6 Summary and Conclusion**

637 The PCMDI has actively developed the PMP with support from the U.S. DOE to improve the understanding
638 of ESMs and to provide systematic and objective ESM evaluation capabilities. With its focus on physical climate, the
639 current evaluation categories enabled in the PMP include seasonal and annual climatology of multiple variables,
640 ENSO, various variability modes in the climate system, MJO, monsoon, cloud feedback and mean state, and simulated
641 precipitation characteristics. The PMP provides quasi-operational ESM evaluation capabilities that can be rapidly
642 deployed to objectively summarize a diverse suite of model behavior with results made publicly available. This can
643 be of value in the assessment of community intercomparisons like CMIP, the evaluation of large ensembles, or the
644 model development process. By documenting objective performance summaries produced by the PMP and making
645 them available via detailed version control, additional research is made possible beyond the baseline model evaluation,
646 model intercomparison, and benchmarking. The outcomes of PMP's calculations applied to the CMIP archive
647 culminate in the PCMDI Simulation Summary (<https://pcmdi.llnl.gov/metrics/>) that has served as a comprehensive
648 data portal for objective model-to-observation comparisons and model-to-model benchmarking and intercomparisons.
649 Special attention is dedicated to the most recent ensemble of models contributing to CMIP6. By offering a diverse and
650 comprehensive suite of evaluation capabilities, the PMP framework equips model developers with quantifiable
651 benchmarks to validate and enhance model performance.

652 We expect that the PMP will continue to play a crucial role in benchmarking ESMs. Improvements in the
653 PMP, along with progress in interconnected MIP community projects, will greatly contribute to advancing the
654 evaluation of ESMs including in connection to the community efforts (e.g., the CMIP Benchmarking Task Team).
655 Enhancements in version control and transparency within obs4MIPs are set to enhance the provenance and
656 reproducibility of PMP results, thereby strengthening the foundation for rigorous and repeatable performance
657 benchmarking. The PMP's collaboration with the CMIP Forcing Task Team, through the Input4MIPs (Durack et al.,
658 2018) and the CMIP6Plus projects, will further expand the utility of performance metrics in identifying problems
659 associated with the forcing dataset and their application and use in reproducing the observed record of historical
660 climate. Furthermore, as ESMs advance towards more operationalized configurations to meet the demands of decision-
661 making processes (Jakob et al., 2023), the PMP holds significant potential to provide interoperable ESM evaluation
662 and benchmarking capabilities to the community.

663 **Appendix A: Table of acronyms**

664

Acronym	Description
AMIP	Atmospheric Model Intercomparison Project
AMO	Atlantic Multi-decadal Oscillation
ARM	Atmospheric Radiation Measurement
ASoP	Analyzing Scales of Precipitation
CBF	Common Basis Function
CDAT	Community Data Analysis Tools
CIESM	Community Integrated Earth System Model
CLIVAR	Climate and Ocean Variability, Predictability and Change
CMEC	Coordinated Model Evaluation Capabilities
CMIP	Coupled Model Intercomparison Project
CMOR	Climate Model Output Rewriter
CVDP	Climate Variability Diagnostics Package
DOE	U.S. Department of Energy
ENSO	El Niño-Southern Oscillation
EOF	Empirical Orthogonal Functions
EOR	East power normalized by Observation

ESGF	Earth System Grid Federation
ESM	Earth System Model
ESMAC Diags	Earth System Model Aerosol–Cloud Diagnostics
ETMoV	Extratropical modes of variability
EWV	East/West power Ratio
GFDL	Geophysical Fluid Dynamics Laboratory
ILAMB	International Land Model Benchmarking
IOMB	International Ocean Model Benchmarking
IPCC	Intergovernmental Panel on Climate Change
IPSL	Institut Pierre-Simon Laplace
ISCCP HGG	International Satellite Cloud Climatology Project H-series Gridded Global
ITCZ	Intertropical Convergence Zone
JSON	JavaScript Object Notation
MAE	Mean Absolute Error
MDTF	Model Diagnostics Task Force
MIPs	Model Intercomparison Projects
MJO	Madden-Julian Oscillation
NAM	Northern Annular Mode

NAO	North Atlantic Oscillation
NASA	National Aeronautics and Space Administration
NCAR	National Center for Atmospheric Research
NetCDF	Network Common Data Form
NH	Northern Hemisphere
NIMS-KMA	National Institute of Meteorological Sciences-Korea Meteorological Administration
NOAA	National Oceanic and Atmospheric Administration
NGO	North Pacific Gyre Oscillation
NPO	North Pacific Oscillation
PCMDI	Program for Climate Model Diagnosis and Intercomparison
PDO	Pacific Decadal Oscillation
PMP	PCMDI Metrics Package
PNA	Pacific North America pattern
RCMES	Regional Climate Model Evaluation System
RMSE	Root-Mean-Square Error
SAM	Southern Annular Mode
SH	Southern Hemisphere
SST	Sea Surface Temperature

TOA	Top of Atmosphere
WCRP	World Climate Research Programme
WGCM	Working Group on Coupled Models
WGNE	Working Group on Numerical Experimentation
xCDAT	Xarray Climate Data Analysis Tools

666 **Code and Data Availability**

667 The source code of the PMP (Lee et al., 2023b) is available as an open-source Python package:
668 https://github.com/PCMDI/pcmdi_metrics (last access: 21 February 2024) with all released versions archived on
669 Zenodo DOI: <https://doi.org/10.5281/zenodo.592790> (last access: 21 February 2024). The online documentation is
670 available at http://pcmdi.github.io/pcmdi_metrics (last access: 21 February 2024). The PMP results database (Lee et
671 al., 2023a) that includes calculated metrics is available on the GitHub repository at
672 https://github.com/PCMDI/pcmdi_metrics_results_archive (last access: 21 February 2024) with versions archived on
673 Zenodo DOI: <https://doi.org/10.5281/zenodo.10181201>. PMP's installation process is streamlined using the *Anaconda*
674 distribution and the *conda-forge* channel (https://anaconda.org/conda-forge/pcmdi_metrics, last access: 21 February
675 2024). The installation instructions are available at http://pcmdi.github.io/pcmdi_metrics/install.html (last access: 21
676 February 2024). The interactive visualizations of the PMP results are available on the PCMDI website at
677 <https://pcmdi.llnl.gov/metrics> (last access: 21 November 2023). The CMIP5 and CMIP6 model outputs and obs4MIPs
678 datasets used in this paper are available via the Earth System Grid Federation at <https://esgf-node.llnl.gov/> (last access:
679 21 February 2024).

680

681 **Author Contributions**

682 All authors contributed to the design and implementation of the research, analysis of the results, and to writing of the
683 manuscript. All authors contributed to the development of codes/metrics in the PMP, its ecosystem tools, and/or the
684 establishment of use cases. JL and PJG led and coordinated the paper with input from all authors.

685

686 **Competing interests**

687 At least one of the coauthors is a member of the editorial board of *Geosic. Model Dev.*. The peer-review process was
688 guided by an independent editor, and the authors also have no other competing interests to declare.

689

690 **Acknowledgment**

691 We acknowledge the World Climate Research Programme, which, through its Working Group on Coupled Modeling,
692 coordinated and promoted CMIP6. We thank the climate modeling groups for producing and making available their
693 model output, the Earth System Grid Federation (ESGF) for archiving the data and providing access, and the multiple
694 funding agencies that support CMIP6 and ESGF. This work is performed under the auspices of the U.S. DOE by
695 Lawrence Livermore National Laboratory (LLNL) under Contract No. DE-AC52-07NA27344. Efforts of JL, PJG,
696 MA, AO, PU, KET, PD, CB, MDZ, LC, and BD were supported by the Regional and Global Model Analysis (RGMA)
697 program of the U.S. Department of Energy (DOE) Office of Science (OS), Biological and Environmental Research
698 (BER) program. MFW was supported by the Director, OS, BER of the U.S. DOE through the RGMA program under
699 Contract No. DE340AC02-05CH11231. AGP was supported by U.S. DOE through BER RGMA through Award
700 Number DE-SC0022070 and via National Science Foundation (NSF) IA 1947282, and by National Center for
701 Atmospheric Research (NCAR), which is a major facility sponsored by the NSF under Cooperative Agreement No.
702 1852977. YYP and EG were supported by the Agence Nationale de la Recherche ARISE project, under Grant ANR-

Deleted: *Geoscientific Model Development*

704 18-CE01-0012, and the Belmont project GOTHAM, under Grant ANR-15-JCLI-0004-01, the European
705 Commission's H2020 Programme "Infrastructure for the European Network for Earth System Modelling Phase 3 (IS-
706 ENES3)" project under Grant Agreement 824084. DK was supported by the New Faculty Startup Fund from Seoul
707 National University and the KMA R&D program (KMI2022-01313). The authors thank Program Manager Renu
708 Joseph of the U.S. DOE for the support and advocacy for the Program for Climate Model Diagnosis and
709 Intercomparison (PCMDI) project and the PMP. We thank Stephen Klein for his leadership for the PCMDI project
710 from 2019 to 2022. We acknowledge contributions from our LLNL colleagues, Lina Muryanto and Zeshawn Shaheen
711 (Now at Google LLC) during the early stage of the PMP, and Sasha Ames, Jeff Painter, Chris Mauzey, and Stephen
712 Po-Chedley for the PCMDI's CMIP database management. The authors also thank Liping Zhang for her comments
713 during GFDL's internal review process.

714

715 **References**

716 Adler, R.F., Sapiano, M. R., Huffman, G.J., Wang, J.J., Gu, G., Bolvin, D., Chiu, L., Schneider, U., Becker, A., Nelkin,
717 E., Xie, P., Ferraro, R., Shin, D.-B.: The Global Precipitation Climatology Project (GPCP) monthly analysis
718 (new version 2.3) and a review of 2017 global precipitation. *Atmosphere*, 9, 138,
719 <https://doi.org/10.3390/atmos9040138>, 2018.

720 Ahn, M.-S., Kim, D. H., Sperber, K. R., Kang, I.-S., Maloney, E. D., Waliser, D. E., and Hendon, H. H.: MJO
721 simulation in CMIP5 climate models: MJO skill metrics and process-oriented diagnosis, *Clim. Dynam.*, 49,
722 4023–4045, <https://doi.org/10.1007/s00382-017-3558-4>, 2017.

723 Ahn, M.-S., Gleckler, P. J., Lee, J., Pendergrass, A. G., and Jakob, C.: Benchmarking Simulated Precipitation
724 Variability Amplitude across Time Scales, *J. Climate*, 35, 3173–3196, [https://doi.org/10.1175/jcli-d-21-](https://doi.org/10.1175/jcli-d-21-0542.1)
725 0542.1, 2022.

726 Ahn, M.-S., Ullrich, P. A., Gleckler, P. J., Lee, J., Ordonez, A. C., and Pendergrass, A. G.: Evaluating precipitation
727 distributions at regional scales: a benchmarking framework and application to CMIP5 and 6 models, *Geosic.*
728 *Model Dev.*, 16, 3927–3951, <https://doi.org/10.5194/gmd-16-3927-2023>, 2023.

729 Arcodia, M., Barnes, E. A., Mayer, K., Lee, J., Ordonez, A., and Ahn, M.-S.: Assessing decadal variability of
730 subseasonal forecasts of opportunity using explainable AI, *Environ. Res.*, 2, 045002,
731 <https://doi.org/10.1088/2752-5295/aced60>, 2023.

732 Ashfaq, M., Rastogi, D., Kitson, J., Abid, M. A., and Kao, S.-C.: Evaluation of CMIP6 GCMs over the CONUS for
733 downscaling studies, *J. Geophys. Res.-Atmos.*, 127, e2022JD036659,
734 <https://doi.org/10.1029/2022JD036659>, 2022.

735 Bayr, T., Wengel, C., Latif, M., Dommenges, D., Lübbecke, J., and Park, W.: Error compensation of ENSO
736 atmospheric feedbacks in climate models and its influence on simulated ENSO dynamics, *Clim. Dynam.*, 53,
737 155–172, <https://doi.org/10.1007/s00382-018-4575-7>, 2019.

738 Biard, J. C. and Kunkel, K. E.: Automated detection of weather fronts using a deep learning neural network, *Adv.*
739 *Stat. Clim. Meteorol. Oceanogr.*, 5, 147–160, <https://doi.org/10.5194/ascmo-5-147-2019>, 2019.

Deleted: Climate Dynamics

Deleted: Journal of Climate

Deleted: Geoscientific Model Development

Deleted: Environmental Research

Deleted: Journal of Geophysical Research: Atmospheres

Deleted: Climate Dynamics

746 Bellenger, H., Guilyardi, E., Leloup, J., Lengaigne, M., and Vialard, J.: ENSO representation in climate models: from
747 CMIP3 to CMIP5, *Clim. Dynam.*, 42, 1999–2018, <https://doi.org/10.1007/s00382-013-1783-z>, 2013.

748 Boucher, O., Servonnat, J., Albright, A. L., Aumont, O., Balkanski, Y., Bastrikov, V., Bekki, S., Bonnet, R., Bony,
749 S., Bopp, L., Braconnot, P., Brockmann, P., Cadule, P., Caubel, A., Cheruy, F., Codron, F., Cozic, A., Cugnet,
750 D., D’Andrea, F., Davini, P., De Lavergne, C., Denvil, S., Deshayes, J., Devilliers, M., Ducharne, A.,
751 Dufresne, J. L., Dupont, E., Ethé, C., Fairhead, L., Falletti, L., Flavoni, S., Foujols, M. A., Gardoll, S.,
752 Gastineau, G., Ghattas, J., Grandpeix, J. Y., Guenet, B., Lionel, E. G., Guilyardi, E., Guimberteau, M.,
753 Hauglustaine, D., Hourdin, F., Idelkadi, A., Joussaume, S., Kageyama, M., Khodri, M., Krinner, G., Lebas,
754 N., Levavasseur, G., Lévy, C., Li, L., Lott, F., Lurton, T., Luyssaert, S., Madec, G., Madeleine, J.-B.,
755 Maignan, F., Marchand, M., Marti, O., Mellul, L., Meurdesoif, Y., Mignot, J., Musat, I., Otlé, C., Peylin, P.,
756 Planton, Y., Polcher, J., Rio, C., Rochetin, N., Rousset, C., Sepulchre, P., Sima, A., Swingedouw, D.,
757 Thiéblemont, R., Traore, A. K., Vancoppenolle, M., Vial, J., Vialard, J., Viovy, N., and Vuichard, N.:
758 Presentation and evaluation of the IPSL-CM6A-LR Climate Model, *J. Adv. Model. Earth Sy.*, 12,
759 <https://doi.org/10.1029/2019ms002010>, 2020.

760 Caldwell, P., Mametjanov, A., Tang, Q., Van Roekel, L., Golaz, J.-C., Lin, W., Bader, D. C., Keen, N. D., Feng, Y.,
761 Jacob, R., Maltrud, M., Roberts, A., Taylor, M. A., Veneziani, M., Wang, H., Wolfe, J. D., Balaguru, K.,
762 Cameron-Smith, P. J., Dong, L., Klein, S. A., Leung, L. R., Li, H., Li, Q., Liu, X., Neale, R., Pinheiro, M.
763 C., Qian, Y., Ullrich, P. A., Xie, S., Yang, Y., Zhang, Y., Zhang, K., and Zhou, T.: The DOE E3SM Coupled
764 Model Version 1: description and results at high resolution, *J. Adv. Model. Earth Sy.*, 11, 4095–4146,
765 <https://doi.org/10.1029/2019ms001870>, 2019.

766 Chen, H.-C., Jin, F.-F., Zhao, S., Wittenberg, A. T., and Xie, S.: ENSO dynamics in the E3SM-1-0, CESM2, and
767 GFDL-CM4 climate models, *J. Climate*, 34, 9365–9384, <https://doi.org/10.1175/JCLI-D-21-0355.1>, 2021.

768 Collier, N., Hoffman, F. M., Lawrence, D. M., Keppel-Aleks, G., Koven, C. D., Riley, W. J., Mu, M., and Randerson,
769 J. T.: The International Land Model Benchmarking (ILAMB) System: Design, theory, and implementation,
770 *J. Adv. Model. Earth Sy.*, 10, 2731–2754, <https://doi.org/10.1029/2018ms001354>, 2018.

771 Covey, C., AchutaRao, K., Cubasch, U., Jones, P., Lambert, S. J., Mann, M., Phillips, T. J., and Taylor, K. E.: An
772 overview of results from the Coupled Model Intercomparison Project, *Global. Planet. Change.*, 37, 103–133,
773 [https://doi.org/10.1016/s0921-8181\(02\)00193-5](https://doi.org/10.1016/s0921-8181(02)00193-5), 2003.

774 Covey, C., Gleckler, P. J., Doutriaux, C., Williams, D. N., Dai, A., Fasullo, J. T., Trenberth, K. E., and Berg, A.:
775 Metrics for the diurnal cycle of precipitation: toward routine benchmarks for climate models, *J. Climate*, 29,
776 4461–4471, <https://doi.org/10.1175/jcli-d-15-0664.1>, 2016.

777 Crockford, D.: The application/json media type for javascript object notation (json) (No. rfc4627), [https://www.rfc-
778 editor.org/rfc/pdf/rfc4627.txt.pdf](https://www.rfc-editor.org/rfc/pdf/rfc4627.txt.pdf) (last access: 5 March 2024), 2006.

779 Crockford, D. and Morningstar, C.: The JSON Data Interchange Syntax, ECMA-404, ECMA International, 2017.

780 Dalelane, C., Winderlich, K., and Walter, A.: Evaluation of global teleconnections in CMIP6 climate projections using
781 complex networks, *Earth Syst. Dynam.*, 14, 17–37, <https://doi.org/10.5194/esd-14-17-2023>, 2023.

Deleted: Climate Dynamics

Deleted: Journal of Advances in Modeling Earth Systems

Deleted: Journal of Advances in Modeling Earth Systems

Deleted: Journal of Climate

Deleted: Journal of Advances in Modeling Earth Systems

Deleted: Global and Planetary Change

Deleted: Journal of Climate

789 Dawson, A.: eofs: A Library for EOF Analysis of Meteorological, Oceanographic, and Climate Data, *J. Open Res.*
790 Software, 4, e14, <https://doi.org/10.5334/jors.122>, 2016.

791 Dec, D. P., Uppala, S. M., Simmons, A. J., Berrisford, P., Poli, P., Kobayashi, S., Andrae, U., Balmaseda, M. A.,
792 Balsamo, G., Bauer, P., Bechtold, P., Beljaars, A. C. M., van de Berg, L., Bidlot, J., Bormann, N., Delsol,
793 C., Dragani, R., Fuentes, M., Geer, A. J., Haimberger, L., Healy, S. B., Hersbach, H., Hólm, E. V., Isaksen,
794 L., Kállberg, P., Köhler, M., Matricardi, M., McNally, A. P., Monge-Sanz, B. M., Morcrette, J.-J., Park, B.-
795 K., Peubey, C., de Rosnay, P., Tavolato, C., Thépaut, J.-N., and Vitart, F.: The ERA-Interim reanalysis:
796 Configuration and performance of the data assimilation system, *Q. J. Roy. Meteor. Soc.*, 137, 553–597,
797 <https://doi.org/10.1002/qj.828>, 2011

798 Deser, C. and Phillips, A. S.: Defining the internal component of Atlantic multidecadal variability in a changing
799 climate, *Geophys. Res. Lett.*, 48, <https://doi.org/10.1029/2021gl095023>, 2021.

800 Doutriaux, C., Nadeau, D., Wittenburg, S., Lipsa, D., Muryanto, L., Chaudhary, A., and Williams, D. N.: CDAT/cdat:
801 CDAT 8.1, Zenodo [Code], <https://doi.org/10.5281/zenodo.2586088>, 2019.

802 Durack, P. J., Taylor, K. E., Eyring, V., Ames, S., Hoang, T., Nadeau, D., Doutriaux, C., Stockhause, M., and Gleckler,
803 P. J.: Toward standardized data sets for climate model experimentation, *Eos, Transactions American*
804 *Geophysical Union*, 99, <https://doi.org/10.1029/2018eo101751>, 2018.

805 Eaton, B., Gregory, J., Drach, B., Taylor, K., Hankin, S., Blower, J., Caron, J., Signell, R., Bentley, P., Rappa, G.,
806 Höck, H., Pamment, A., Juckes, M., Raspaud, M., Horne, R., Whiteaker, T., Blodgett, D., Zender, C., Lee,
807 D., Hassell, D., Snow, A. D., Kölling, T., Allured, D., Jelenak, A., Soerensen, A. M., Gaultier, L., Herlédan,
808 S.: NetCDF Climate and Forecast (CF) Meta-data Conventions V1.10, available at:
809 <http://cfconventions.org/Data/cf-conventions/cf-conventions-1.10/cf-conventions.html> (last access: 6
810 November 2023), 2022.

811 Eyring, V., Righi, M., Lauer, A., Evaldsson, M., Wenzel, S., Jones, C., Anav, A., Andrews, O., Cionni, I., Davin, E.
812 L., Deser, C., Ehbrecht, C., Friedlingstein, P., Gleckler, P. J., Gottschaldt, K.-D., Hagemann, S., Juckes, M.,
813 Kindermann, S., Krasting, J. P., Kunert, D., Levine, R. C., Loew, A., Mäkelä, J., Martin, G., Mason, E.,
814 Phillips, A. S., Read, S., Rio, C., Roehrig, R., Senftleben, D., Sterl, A., Van Ulft, L. H., Walton, J., Wang,
815 S., and Williams, K. D.: ESMValTool (v1.0) – a community diagnostic and performance metrics tool for
816 routine evaluation of Earth system models in CMIP, *Geosic. Model Dev.*, 9, 1747–1802,
817 <https://doi.org/10.5194/gmd-9-1747-2016>, 2016a.

818 Eyring, V., Bony, S., Meehl, G. A., A. C., Senior, Stevens, B., Stouffer, R. J., and Taylor, K. E.: Overview of the
819 Coupled Model Intercomparison Project Phase 6 (CMIP6) experimental design and organization, *Geosic.*
820 *Model Dev.*, 9, 1937–1958, <https://doi.org/10.5194/gmd-9-1937-2016>, 2016b.

821 Eyring, V., Cox, P. M., Flato, G. M., Gleckler, P. J., Abramowitz, G., Caldwell, P., Collins, W. D., Gier, B. K., Hall,
822 A., Hoffman, F. M., Hurtt, G. C., Jahn, A., Jones, C. D., Klein, S. A., Krasting, J. P., Kwiatkowski, L.,
823 Lorenz, R., Maloney, E. D., Meehl, G. A., Pendergrass, A. G., Pincus, R., Ruane, A. C., Russell, J. L.,
824 Sanderson, B. M., Santer, B. D., Sherwood, S. C., Simpson, I. R., Stouffer, R. J., and Williamson, M. S.:

Deleted: Journal

Deleted: of

Deleted: Research

Deleted: (JORS)

Deleted: Quarterly Journal of the Royal Meteorological Society...

Deleted: Geophysical Research Letters

Deleted: Geoscientific Model Development

Deleted: Geoscientific Model Development

834 Taking climate model evaluation to the next level, [Nat. Clim. Change](#), 9, 102–110,
835 <https://doi.org/10.1038/s41558-018-0355-y>, 2019.

836 Eyring, V., Bock, L., Lauer, A., Righi, M., Schlund, M., Andela, B., Arnone, E., Bellprat, O., Brötz, B., Caron, L.-P.,
837 Carvalhais, N., Cionni, I., Cortesi, N., Crezee, B., Davin, E. L., Davini, P., Debeire, K., De Mora, L., Deser,
838 C., Docquier, D., Earnshaw, P., Ehbrecht, C., Gier, B. K., Gonzalez-Reviriego, N., Goodman, P. J.,
839 Hagemann, S., Hardiman, S. C., Hassler, B., Hunter, A., Kadow, C., Kindermann, S., Koirala, S., Koldunov,
840 N., Lejeune, Q., Lembo, V., Lovato, T., Lucarini, V., Massonnet, F., Müller, B., Pandde, A., Pérez-Zanón,
841 N., Phillips, A. S., Predoi, V., Russell, J. L., Sellar, A., Serva, F., Stacke, T., Swaminathan, R., Torralba, V.,
842 Vegas-Regidor, J., Von Hardenberg, J., Weigel, K., and Zimmermann, K.: Earth System Model Evaluation
843 Tool (ESMValTool) v2.0 – an extended set of large-scale diagnostics for quasi-operational and
844 comprehensive evaluation of Earth system models in CMIP, [Geosic. Model Dev.](#), 13, 3383–3438,
845 <https://doi.org/10.5194/gmd-13-3383-2020>, 2020.

846 Eyring, V., Gillett, N.P., Achuta Rao, K.M., Barimalala, R., Barreiro Parrillo, M., Bellouin, N., Cassou, C., Durack,
847 P.J., Kosaka, Y., McGregor, S. and Min, S., Morgenstern, O., and Sun, Y.: Human Influence on the Climate
848 System. In *Climate Change 2021: The Physical Science Basis. Contribution of Working Group I to the Sixth*
849 *Assessment Report of the Intergovernmental Panel on Climate Change*. 105, 423-552,
850 <https://doi.org/10.1017/9781009157896.005>, 2021.

851 Fasullo, J. T.: Evaluating simulated climate patterns from the CMIP archives using satellite and reanalysis datasets
852 using the Climate Model Assessment Tool (CMATv1), [Geosic. Model Dev.](#), 13, 3627–3642,
853 <https://doi.org/10.5194/gmd-13-3627-2020>, 2020.

854 Fasullo, J. T., Phillips, A. S., and Deser, C.: Evaluation of leading modes of climate variability in the CMIP archives,
855 [J. Climate](#), 33, 5527–5545, <https://doi.org/10.1175/jcli-d-19-1024.1>, 2020.

856 Ferraro, R., Waliser, D. E., Gleckler, P. J., Taylor, K. E., and Eyring, V.: Evolving OBS4MIPS to support Phase 6 of
857 the Coupled Model Intercomparison Project (CMIP6), [B. Am. Meteorol. Soc.](#), [https://doi.org/10.1175/bams-](https://doi.org/10.1175/bams-d-14-00216.1)
858 [d-14-00216.1](https://doi.org/10.1175/bams-d-14-00216.1), 2015.

859 Flato, G., Marotzke, J., Abiodun, B., Braconnot, P., Chou, S.C., Collins, W., Cox, P., Driouech, F., Emori, S., Eyring,
860 V. and Forest, C.: Evaluation of climate models. In *Climate change 2013: the physical science basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change* (pp. 741-866). Cambridge University Press. 2014.

863 Fu, W., Moore, J. K., Primeau, F., Collier, N., Ogunro, O. O., Hoffman, F. M. and Randerson, J. T.: Evaluation of
864 ocean biogeochemistry and carbon cycling in CMIP earth system models with the international ocean model
865 benchmarking (IOMB) software System. [J. Geophys. Res.-Oceans](#), 127, e2022JC018965,
866 <https://doi.org/10.1029/2022JC018965>, 2022.

867 Gates, W.L.: AN AMS continuing series: Global CHANGE–AMIP: The Atmospheric Model Intercomparison Project,
868 [B. Am. Meteorol. Soc.](#), 73, 1962-1970, 1992.

869 Gates, W.L., Henderson-Sellers, A., Boer, G.J., Folland, C.K., Kitoh, A., McAvaney, B.J., Semazzi, F., Smith, N.,
870 Weaver, A.J. and Zeng, Q.C.: Climate models—evaluation. *Climate Change* 1: 229-284, 1995.

Deleted: Nature Climate Change

Deleted: Geoscientific Model Development

Deleted: Geoscientific Model Development

Deleted: Journal of Climate

Deleted: Bulletin of the American Meteorological Society

Deleted: Journal of Geophysical Research: Oceans

Deleted: Bulletin of the American Meteorological Society

Deleted: change

879 Gates, W.L., Boyle, J.S., Covey, C., Dease, C.G., Doutriaux, C.M., Drach, R.S., Fiorino, M., Gleckler, P.J., Hnilo,
880 J.J., Marlais, S.M. and Phillips, T.J.: An overview of the results of the Atmospheric Model Intercomparison
881 Project (AMIP I). *B. Am. Meteorol. Soc.*, 80, 29-56, 1999.

882 Gleckler, P. J., Taylor, K. E., and Doutriaux, C.: Performance metrics for climate models, *J. Geophys. Res.*, 113,
883 <https://doi.org/10.1029/2007jd008972>, 2008.

884 Gleckler, P. J., Ferraro, R., and Waliser, D. E.: Improving use of satellite data in evaluating climate models, *Eos*,
885 *Transactions American Geophysical Union*, 92, 172, <https://doi.org/10.1029/2011eo200005>, 2011.

886 Gleckler, P. J., Doutriaux, C., Durack, P. J., Taylor, K. E., Zhang, Y., Williams, D. N., Mason, E., and Servonnat, J.:
887 A more powerful reality test for climate models, *Eos, Transactions American Geophysical Union*, 97,
888 <https://doi.org/10.1029/2016eo051663>, 2016.

889 Golaz, J.-C., Caldwell, P., Van Roekel, L., Petersen, M. R., Tang, Q., Wolfe, J. D., Abeshu, G. W., Anantharaj, V.,
890 Asay-Davis, X., Bader, D. C., Baldwin, S., Bisht, G., Bogenschutz, P., Branstetter, M. L., Brunke, M. A.,
891 Brus, S., Burrows, S. M., Cameron-Smith, P. J., Donahue, A. S., Deakin, M., Easter, R. C., Evans, K. J.,
892 Feng, Y., Flanner, M., Foucar, J. G., Fyke, J., Griffin, B. M., Hannay, C., Harrop, B. E., Hoffman, M. J.,
893 Hunke, E., Jacob, R., Jacobsen, D. W., Jeffery, N., Jones, P. W., Keen, N. D., Klein, S. A., Larson, V. E.,
894 Leung, L. R., Li, H. Y., Lin, W., Lipscomb, W. H., Lun, P., Mahajan, S., Maltrud, M., Mامتjanov, A.,
895 McClean, J. L., McCoy, R., Neale, R., Price, S., Qian, Y., Rasch, P. J., Eyre, J. E. J. R., Riley, W. J., Ringler,
896 T. D., Roberts, A., Roesler, E. L., Salinger, A. G., Shaheen, Z., Shi, X., Singh, B., Tang, J., Taylor, M. A.,
897 Thornton, P. E., Turner, A. K., Veneziani, M., Wan, H., Wang, H., Wang, S., Williams, D. N., Wolfram, P.
898 J., Worley, P. H., Xie, S., Yang, Y., Yoon, J., Zelinka, M. D., Zender, C. S., Zeng, X., Zhang, C., Zhang, K.,
899 Zhang, Y., Zheng, X., Zhou, T., and Zhu, Q.: The DOE E3SM Coupled Model Version 1: Overview and
900 evaluation at standard resolution, *J. Adv. Model. Earth Sy.*, 11, 2089–2129,
901 <https://doi.org/10.1029/2018ms001603>, 2019.

902 Goldenson, N., Leung, L. R., Mearns, L. O., Pierce, D. W., Reed, K. A., Simpson, I. R., Ullrich, P., Krantz, W., Hall,
903 A., Jones, A. and Rahimi, S.: Use-Inspired, Process-Oriented GCM Selection: Prioritizing Models for
904 Regional Dynamical Downscaling, *B. Am. Meteorol. Soc.*, E1619–E1629, [https://doi.org/10.1175/BAMS-](https://doi.org/10.1175/BAMS-D-23-0100.1)
905 [D-23-0100.1](https://doi.org/10.1175/BAMS-D-23-0100.1), 2023.

906 Guilyardi, E., Wittenberg, A., Fedorov, A., Collins, M., Wang, C., Capotondi, A., Van Oldenborgh, G.J. and
907 Stockdale, T.: Understanding El Niño in ocean–atmosphere general circulation models: Progress and
908 challenges, *B. Am. Meteorol. Soc.*, 90, 325-340, <https://doi.org/10.1175/2008BAMS2387.1>, 2009.

909 Guilyardi E., Capotondi, A., Lengaigne, M., Thual, S., Wittenberg, A. T.: ENSO modelling: history, progress and
910 challenges, in: *El Niño in a changing climate*, edited by: McPhaden, M. J., Santoso, A., Cai, W., AGU
911 monograph, ISBN: 9781119548164, <https://doi.org/10.1002/9781119548164.ch9>, 2020.

912 Gutowski Jr., W. J., Giorgi, F., Timbal, B., Frigon, A., Jacob, D., Kang, H.-S., Raghavan, K., Lee, B., Lennard, C.,
913 Nikulin, G., O'Rourke, E., Rixen, M., Solman, S., Stephenson, T., and Tangang, F.: WCRP Coordinated
914 Regional Downscaling EXperiment (CORDEX): a diagnostic MIP for CMIP6, *Geosci. Model Dev.*, 9, 4087–
915 4095, <https://doi.org/10.5194/gmd-9-4087-2016>, 2016.

Deleted: Bulletin of the American Meteorological Society

Deleted: Journal of Geophysical Research

Deleted: Journal of Advances in Modeling Earth Systems

Deleted: Bulletin of the American Meteorological Society

Deleted: Bulletin of the American Meteorological Society

921 Haarsma, R. J., Roberts, M., Vidale, P. L., A. C., Senior, Bellucci, A., Bao, Q., Chang, P., Corti, S., Fučkar, N. S.,
 922 Guemas, V., Von Hardenberg, J., Hazeleger, W., Kodama, C., Koenigk, T., Leung, L. R., Lu, J., Luo, J.,
 923 Mao, J., Mizielinski, M. S., Mizuta, R., Nobre, P., Satoh, M., Scoccimarro, E., Semmler, T., Small, R. J., and
 924 Von Storch, J. S.: High Resolution Model Intercomparison Project (HiGHRESMIP v1.0) for CMIP6, *Geosic.*
 925 *Model Dev.*, 9, 4185–4208, <https://doi.org/10.5194/gmd-9-4185-2016>, 2016.

926 Hannah, W. M., Bradley, A. M., Guba, O., Tang, Q., Golaz, J.-C., and Wolfe, J. D.: Separating physics and dynamics
 927 grids for improved computational efficiency in spectral element Earth system models, *J. Adv. Model. Earth*
 928 *Sy.*, 13, <https://doi.org/10.1029/2020ms002419>, 2021.

929 Hassan, K. A., Rönnerberg, N., Forsell, C., Cooper, M. and Johansson, J.: A study on 2D and 3D parallel coordinates
 930 for pattern identification in temporal multivariate data, in: 2019 23rd International Conference Information
 931 Visualisation (IV), 145-150, <https://doi.org/10.1109/IV.2019.00033>, 2019.

932 Hassell, D., Gregory, J. M., Blower, J., Lawrence, B., and Taylor, K. E.: A data model of the Climate and Forecast
 933 metadata conventions (CF-1.6) with a software implementation (cf-python v2.1), *Geosic. Model Dev.*, 10,
 934 4619–4646, <https://doi.org/10.5194/gmd-10-4619-2017>, 2017.

935 Hausfather, Z., Marvel, K., Schmidt, G. A., Nielsen-Gammon, J. W. and Zelinka, M.: Climate simulations: Recognize
 936 the ‘hot model’ problem, *Nature*, 605, 26-29, <https://doi.org/10.1038/d41586-022-01192-2>, 2022.

937 Held, I. M., Guo, H., Adcroft, A., Dunne, J. P., Horowitz, L. W., Krasting, J., Shevliakova, E., Winton, M., Zhao, M.,
 938 Bushuk, M., Wittenberg, A. T., and coauthors: Structure and performance of GFDL’s CM4. 0 climate model,
 939 *J. Adv. Model. Earth Sy.*, 11, 3691-3727, <https://doi.org/10.1029/2019MS001829>, 2019.

940 Hendon, H. H., Zhang, C., and Glick, J. D.: Interannual Variation of the Madden–Julian Oscillation during Austral
 941 Summer, *J. Climate*, 12, 2538–2550, 1999.

942 Herger, N., Abramowitz, G., Knutti, R., Angéllil, O., Lehmann, K., and Sanderson, B. M.: Selecting a climate model
 943 subset to optimise key ensemble properties, *Earth System Dynamics Discussions*, 9, 135–151,
 944 <https://doi.org/10.5194/esd-9-135-2018>, 2018.

945 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., Nicolas, J., Peubey, C., Radu, R.,
 946 Schepers, D. and coauthors: The ERA5 global reanalysis. *Q. J. Roy. Meteor. Soc.*, 146, 1999-2049,
 947 <https://doi.org/10.1002/qj.3803>, 2020.

948 Hintze, J. L., and Nelson, R. D.: Violin plots: A box plot-density trace synergism, *The American Statistician*, 52, 181–
 949 184, <https://doi.org/10.1080/00031305.1998.10480559>, 1998.

950 Hoyer, S. and Hamman, J.: xarray: N-D labeled Arrays and Datasets in Python, *J. Open Res. Software*, 5, 10,
 951 <https://doi.org/10.5334/jors.148>, 2017.

952 Huffman, G. J., Adler, R. F., Morrissey, M. M., Bolvin, D. T., Curtis, S., Joyce, R., McGavock, B. and Susskind, J.:
 953 Global precipitation at one-degree daily resolution from multisatellite observations, *J. Hydrometeorol.*, 2, 36-
 954 50, 2001.

955 Huffman, G. J., Bolvin, D. T., Braithwaite, D., Hsu, K., Joyce, R., Kidd, C., Nelkin, E. J., Sorooshian, S., Tan, J., and
 956 Xie, P.: NASA global precipitation measurement (GPM) integrated multi-satellite retrievals for GPM
 957 (IMERG). Algorithm theoretical basis document (ATBD) version, 4, p.30., 2015.

Deleted: Geoscientific Model Development

Deleted: Journal of Advances in Modeling Earth Systems

Deleted: Geoscientific Model Development

Deleted: Journal of Advances in Modeling Earth Systems

Deleted: Journal of Climate

Deleted: Quarterly Journal of the Royal Meteorological Society...

Deleted: ¶

Formatted: Indent: Hanging: 0.5"

Deleted: Journal

Deleted: of

Deleted: Research

Deleted: Journal of hydrometeorology

970 Inselberg, A.: Multidimensional detective, in: Proceedings of IEEE Symposium on Information Visualization, 100–
971 107, <https://doi.org/10.1109/INFVIS.1997.636793>, 1997.

972 Inselberg, A.: Parallel Coordinates: Visualization, Exploration and Classification of High-Dimensional Data, in:
973 Handbook of Data Visualization, edited by Chen, C., Härdle, W., and Unwin, A., Springer, Berlin,
974 Heidelberg, Germany, 643–680, https://doi.org/10.1007/978-3-540-33037-0_25, 2008.

975 Inselberg, A.: Parallel Coordinates, in: Encyclopedia of Database Systems. Springer, edited by Liu, L., and Özsu, M.
976 T., Springer, New York, NY, U.S.A., https://doi.org/10.1007/978-1-4899-7993-3_262-2, 2016.

977 Johansson, J. and Forsell, C.: Evaluation of parallel coordinates: Overview, categorization and guidelines for future
978 research, *IEEE T. Vis. Comput. G. R.*, 22, 579–588, <https://doi.org/10.1109/TVCG.2015.2466992>, 2016.

979 Jakob, C., Gettelman, A. and Pitman, A.: The need to operationalize climate modelling, *Nat. Clim. Chang.* 13, 1158–
980 1160, <https://doi.org/10.1038/s41558-023-01849-4>, 2023.

981 Joyce, R. J., Janowiak, J. E., Arkin, P. A. and Xie, P.: CMORPH: A method that produces global precipitation
982 estimates from passive microwave and infrared data at high spatial and temporal resolution, *J.*
983 *Hydrometeorol.*, 5, 487–503, 2004.

984 Kang, D., Kim, D. H., Ahn, M.-S., Neale, R., Lee, J., and Gleckler, P. J.: The role of the mean state on MJO simulation
985 in CESM2 ensemble simulation, *Geophys. Res. Lett.*, 47, <https://doi.org/10.1029/2020gl089824>, 2020.

986 Kim, D., Sperber, K. R., Stern, W., Waliser, D. E., Kang, I. S., Maloney, E. D., Wang, W., Weickmann, K. M.,
987 Benedict, J. J., Khairoutdinov, M., Lee, M.-I., Neale, R., Suarez, M. J., Thayer-Calder, K., and Zhang, G.:
988 Application of MJO simulation diagnostics to climate models, *J. Climate*, 22, 6413–6436,
989 <https://doi.org/10.1175/2009jcli3063.1>, 2009.

990 Kim, H., Caron, J. M., Richter, J. H. and Simpson, I. R.: The lack of QBO-MJO connection in CMIP6 models,
991 *Geophys. Res. Lett.*, 47, e2020GL087295, <https://doi.org/10.1029/2020GL087295>, 2020.

992 Klein, S. A., Zhang, Y., Zelinka, M. D., Pincus, R., Boyle, J., and Gleckler, P. J.: Are climate model simulations of
993 clouds improving? An evaluation using the ISCCP simulator, *J. Geophys. Res.-Atmos.*, 118, 1329–1342,
994 <https://doi.org/10.1002/jgrd.50141>, 2013.

995 Klingaman, N. P., Martin, G., and Moise, A.: ASoP (v1.0): a set of methods for analyzing scales of precipitation in
996 general circulation models, *Geosic. Model Dev.*, 10, 57–83, <https://doi.org/10.5194/gmd-10-57-2017>, 2017.

997 Knutti, R.: The end of model democracy? *Climatic Change*, 102, 395–404, <https://doi.org/10.1007/s10584-010-9800-2>, 2010.

998

999 Knutti, R., Sedláček, J., Sanderson, B. M., Lorenz, R., Fischer, E. M., and Eyring, V.: A climate model projection
1000 weighting scheme accounting for performance and interdependence, *Geophys. Res. Lett.*,
1001 <https://doi.org/10.1002/2016gl072012>, 2017.

1002 Knutti, R., Abramowitz, G., Collins, M., Eyring, V., Gleckler, P. J., Hewitson, B., and Mearns, L.: Good Practice
1003 Guidance Paper on Assessing and Combining Multi Model Climate Projections, in: Meeting Report of the
1004 Intergovernmental Panel on Climate Change Expert Meeting on Assessing and Combining Multi Model
1005 Climate Projections, edited by Stocker, T. F., Qin, D., Plattner, G.-K., Tignor, M., and Midgley, P. M., IPCC
1006 Working Group I Technical Support Unit, University of Bern, Bern, Switzerland, 2010.

Deleted: IEEE Transactions on Visualization and Computer Graphics...

Deleted: Journal of hydrometeorology

Deleted: Geophysical Research Letters

Deleted: Journal of Climate

Deleted: Geophysical Research Letters

Deleted: Journal of Geophysical Research: Atmospheres

Deleted: Geoscientific Model Development

Deleted: Geophysical Research Letters

1016 Labe, Z. M. and Barnes, E. A.: Comparison of Climate Model Large Ensembles With Observations in the Arctic Using
1017 Simple Neural Networks, *Earth and Space Science*, e2022EA002348,
1018 <https://doi.org/10.1029/2022EA002348>, 2022.

1019 Lambert, S. J. and Boer, G. J.: CMIP1 evaluation and intercomparison of coupled climate models, *Clim. Dynam.*, 17,
1020 83-106, <https://doi.org/10.1007/PL00013736>, 2001.

1021 Lee, H., Goodman, A., McGibbney, L. J., Waliser, D. E., Kim, J., Loikith, P. C., Gibson, P. B., and Massoud, E.:
1022 Regional Climate Model Evaluation System powered by Apache Open Climate Workbench v1.3.0: an
1023 enabling tool for facilitating regional climate studies, *Geosic. Model Dev.*, 11, 4435-4449,
1024 <https://doi.org/10.5194/gmd-11-4435-2018>, 2018a.

1025 Lee, J., Ahn, M.-S., Ordonez, A., Gleckler, P., and Ullrich, P.: PCMDI/pcmdi_metrics_results_archive, Zenodo [data],
1026 <https://doi.org/10.5281/zenodo.10181201>, 2023a.

1027 Lee, J., Gleckler, P., Ordonez, A., Ahn, M.-S., Ullrich, P., Tom, V., Jason, B., Charles, D., Durack, P., Shaheen, Z.,
1028 Muryanto, L., Painter, J., and Krasting, J.: PCMDI/pcmdi_metrics: PMP Version 3.1.1, Zenodo [code],
1029 <https://doi.org/10.5281/zenodo.592790>, 2023b.

1030 Lee, J., Gleckler, P., Sperber, K., Doutriaux C., and Williams, D.: High-dimensional Data Visualization for Climate
1031 Model Intercomparison: Application of the Circular Plot, in: Proceedings of the 8th International Workshop
1032 on Climate Informatics: CI 2018. NCAR Technical Note NCAR/TN-550+PROC, 12-14,
1033 <http://dx.doi.org/10.5065/D6BZ64XQ>, 2018b.

1034 Lee, J., Planton, Y., Gleckler, P. J., Sperber, K. R., Guilyardi, E., Wittenberg, A. T., McPhaden, M. J., and Pallotta,
1035 G.: Robust evaluation of ENSO in climate models: How many ensemble members are needed?, *Geophys.*
1036 *Res. Lett.*, 48, <https://doi.org/10.1029/2021gl095041>, 2021a.

1037 Lee, J., Sperber, K. R., Gleckler, P. J., Bonfils, C., and Taylor, K. E.: Quantifying the agreement between observed
1038 and simulated extratropical modes of interannual variability, *Clim. Dynam.*, 52, 4057-4089,
1039 <https://doi.org/10.1007/s00382-018-4355-4>, 2019a.

1040 Lee, J., Sperber, K. R., Gleckler, P. J., Taylor, K. E., and Bonfils, C.: Benchmarking performance changes in the
1041 simulation of extratropical modes of variability across CMIP generations, *J. Climate*, 1-70,
1042 <https://doi.org/10.1175/jcli-d-20-0832.1>, 2021b.

1043 Lee, J., Xue, Y., De Sales, F., Diallo, I., Marx, L., Ek, M., Sperber, K. R., and Gleckler, P. J.: Evaluation of multi-
1044 decadal UCLA-CFSv2 simulation and impact of interactive atmospheric-ocean feedback on global and
1045 regional variability, *Clim. Dynam.*, 52, 3683-3707, <https://doi.org/10.1007/s00382-018-4351-8>, 2019b.

1046 Leung, L. R., Boos, W. R., Catto, J. L., DeMott, C. A., Martin, G. M., Neelin, J. D., O'Brien, T. A., Xie, S., Feng, Z.,
1047 Klingaman, N. P. Kuo, Y.-H., Lee, R. W., Martinez-Villalobos, C., Vishnu S., Priestley, M. D. K., Tao, C.,
1048 and Zhou, Y.: Exploratory precipitation metrics: Spatiotemporal characteristics, process-oriented, and
1049 phenomena-based, *J. Climate*, 35, <https://doi.org/10.1175/JCLI-D-21-0590.1>, 3659-3686, 2022.

1050 Lin, J.-P., Kiladis, G. N., Mapes, B. E., Weickmann, K. M., Sperber, K. R., Lin, W., Wheeler, M. C., Schubert, S. D.,
1051 Del Genio, A. D., Donner, L. J., Emori, S., Guérémy, J.-F., Hourdin, F., Rasch, P. J., Roeckner, E., and

Deleted: Climate Dynamics

Deleted: Geoscientific Model Development

Deleted: Geophysical Research Letters

Deleted: Climate Dynamics

Deleted: Journal of Climate

Deleted: Climate Dynamics

Deleted: Journal of Climate

1059 Scinocca, J.: Tropical intraseasonal variability in 14 IPCC AR4 climate Models. Part I: Convective Signals,
1060 [J. Climate](https://doi.org/10.1175/jcli3735.1), 19, 2665–2690, <https://doi.org/10.1175/jcli3735.1>, 2006.

1061 Lin, Y., Huang, X., Liang, Y., Qin, Y., Xu, S., Huang, W., Xu, F., Liu, L., Wang, Y., Peng, Y. and Wang, L.:
1062 Community integrated earth system model (CIesm): Description and evaluation, [J. Adv. Model. Earth Sy.](https://doi.org/10.1029/2019ms002036),
1063 12, <https://doi.org/10.1029/2019ms002036>, 2020.

1064 Loeb, N. G., Doelling, D. R., Wang, H., Su, W., Nguyen, C., Corbett, J. G., Liang, L., Mitrescu, C., Rose, F. G., and
1065 Seiji, K.: Clouds and the Earth's Radiant Energy System (CERES) Energy Balanced and Filled (EBAF) top-
1066 of-atmosphere (TOA) Edition-4.0 data product, *International Journal of Climatology*, 31, 895–918,
1067 <https://doi.org/10.1175/JCLI-D-17-0208.1>, 2018.

1068 Longmate, J. M., Risser, M. D. and Feldman, D. R.: Prioritizing the selection of CMIP6 model ensemble members for
1069 downscaling projections of CONUS temperature and precipitation, *Clim. Dyn.*, 61, 5171–5197,
1070 <https://doi.org/10.1007/s00382-023-06846-z>, 2023.

1071 Lu, L., Wang, W. and Tan, Z.: Double-arc parallel coordinates and its axes re-ordering methods, *Mobile Networks*
1072 *and Applications*, 25, 1376-1391, <https://doi.org/10.1007/s11036-019-01455-9>, 2020.

1073 Madden, R. A. and Julian, P.: Detection of a 40–50 day oscillation in the zonal wind in the Tropical Pacific, [J. Atmos.](https://doi.org/10.1175/1520-0469(1971)028)
1074 [Sci.](https://doi.org/10.1175/1520-0469(1971)028), 28, 702–708, [https://doi.org/10.1175/1520-0469\(1971\)028](https://doi.org/10.1175/1520-0469(1971)028), 1971.

1075 Madden, R. A. and Julian, P.: Description of Global-Scale Circulation Cells in the Tropics with a 40–50 Day Period,
1076 [J. Atmos. Sci.](https://doi.org/10.1175/1520-0469(1972)029), 29, 1109–1123, [https://doi.org/10.1175/1520-0469\(1972\)029](https://doi.org/10.1175/1520-0469(1972)029), 1972.

1077 Madden, R. A. and Julian, P.: Observations of the 40–50-Day Tropical Oscillation—A Review, [Mon. Weather Rev.](https://doi.org/10.1175/1520-0493(1994)122),
1078 122, 814–837, [https://doi.org/10.1175/1520-0493\(1994\)122](https://doi.org/10.1175/1520-0493(1994)122), 1994.

1079 Martin, G. M., Klingaman, N. P., and Moise, A. F.: Connecting spatial and temporal scales of tropical precipitation in
1080 observations and the MetUM-GA6, [Geosic. Model Dev.](https://doi.org/10.5194/gmd-10-105-2017), 10, 105–126, [https://doi.org/10.5194/gmd-10-105-](https://doi.org/10.5194/gmd-10-105-2017)
1081 2017, 2017.

1082 Maloney, E. D., Gettelman, A., Ming, Y., Neelin, J. D., Barrie, D., Mariotti, A., Chen, C., Coleman, D., Kuo, Y. H.,
1083 Singh, B., Annamalai, H., Berg, A., Booth, J. F., Camargo, S. J., Dai, A., Gonzalez, A., Hafner, J., Jiang, X.,
1084 Jing, X., Kim, D. H., Kumar, A., Moon, Y., Naud, C. M., Sobel, A. H., Suzuki, K., Wang, F., Wang, J., Wing,
1085 A. A., Xu, X., and Zhao, M.: Process-Oriented evaluation of climate and weather forecasting models, [B. Am.](https://doi.org/10.1175/bams-d-18-0042.1)
1086 [Meteorol. Soc.](https://doi.org/10.1175/bams-d-18-0042.1), 100, 1665–1686, <https://doi.org/10.1175/bams-d-18-0042.1>, 2019.

1087 McAvaney, B.J., Covey, C., Joussaume, S., Kattsov, V., Kitoh, A., Ogana, W., Pitman, A.J., Weaver, A.J., Wood,
1088 R.A. and Zhao, Z.C.: Model evaluation. In *Climate Change 2001: The scientific basis. Contribution of WG1*
1089 *to the Third Assessment Report of the IPCC (TAR)* 471-523, Cambridge University Press, 2001.

1090 McPhaden, M. J., Zebiak, S. E., and Glantz, M. H.: ENSO as an integrating concept in Earth Science, *Science*, 314,
1091 1740–1745, <https://doi.org/10.1126/science.1132588>, 2006.

1092 McPhaden, M. J., Santoso, A., Cai, W. (Eds.): *El Niño Southern oscillation in a changing climate*, American
1093 Geophysical Union, USA, 528 pp., ISBN:9781119548126, <https://doi.org/10.1002/9781119548164>, 2020.

Deleted: Journal of Climate

Deleted: Journal of Advances in Modeling Earth Systems

Deleted: Journal of the Atmospheric Sciences

Deleted: Journal of the Atmospheric Sciences

Deleted: Monthly Weather Review

Deleted: Geoscientific Model Development

Deleted: Bulletin of the American Meteorological Society

1101 Mears, C. A., Smith, D. K., Ricciardulli, L., Wang, J., Huelsing, H., & Wentz, F. J.: Construction and uncertainty
1102 estimation of a satellite-derived total precipitable water data record over the world's oceans, *Earth and Space*
1103 *Science*, 5, 197–210, <https://doi.org/10.1002/2018EA000363>, 2018.

1104 Meehl, G. A., Boer, G. J., Covey, C., Latif, M., and Stouffer, R. J.: The Coupled Model Intercomparison Project
1105 (CMIP), *B. Am. Meteorol. Soc.*, 81, 313–318, 2000.

1106 Meehl, G. A., Boer, G. J., Covey, C., Latif, M., and Stouffer, R. J.: Intercomparison makes for a better climate model,
1107 *Eos, Transactions American Geophysical Union*, 78, 445, <https://doi.org/10.1029/97eo00276>, 1997.

1108 Meehl, G. A., Covey, C., Delworth, T. L., Latif, M., McAvaney, B. J., Mitchell, J. F. B., Stouffer, R. J., and Taylor,
1109 K. E.: THE WCRP CMIP3 Multimodel Dataset: A new era in climate change research, *B. Am. Meteorol.*
1110 *Soc.*, 88, 1383–1394, <https://doi.org/10.1175/bams-88-9-1383>, 2007.

1111 Merrifield, A., Brunner, L., Lorenz, R., Humphrey, V., and Knutti, R.: Climate model Selection by Independence,
1112 Performance, and Spread (ClimSIPS v1.0.1) for regional applications, *Geosci. Model Dev.*, 16, 4715–4747,
1113 <https://doi.org/10.5194/gmd-16-4715-2023>, 2023.

1114 Neelin, J. D., Krasting, J. P., Radhakrishnan, A., Liptak, J., Jackson, T. J., Ming, Y., Dong, W., Gettelman, A.,
1115 Coleman, D., Maloney, E. D., Wing, A. A., Kuo, Y. H., Ahmed, F., Ullrich, P. A., Bitz, C. M., Neale, R.,
1116 Ordonez, A., and Maroon, E.: Process-oriented diagnostics: principles, practice, community development
1117 and common standards, *B. Am. Meteorol. Soc.*, <https://doi.org/10.1175/bams-d-21-0268.1>, 2023.

1118 Nowack, P., Runge, J., Eyring, V., and Haigh, J. D.: Causal networks for climate model evaluation and constrained
1119 projections, *Nat. Commun.*, 11, <https://doi.org/10.1038/s41467-020-15195-y>, 2020.

1120 Orbe, C., Van Roekel, L., Adames, Á. F., Dezfúli, A., Fasullo, J. T., Gleckler, P. J., Lee, J., Li, W., Nazarenko, L.,
1121 Schmidt, G. A., Sperber, K. R., and Zhao, M.: Representation of modes of variability in six U.S. climate
1122 models, *J. Climate*, 33, 7591–7617, <https://doi.org/10.1175/jcli-d-19-0956.1>, 2020.

1123 Ordonez, A. C., Klingaman, N. P., and Martin, G.: Analysing scales of precipitation, OSTI OAI (U.S. Department of
1124 Energy Office of Scientific and Technical Information), <https://doi.org/10.11578/dc.20211029.5>, 2021.

1125 Papalexiou, S. M., Rajulapati, C. R., Clark, M. P., and Lehner, F.: Robustness of CMIP6 historical global mean
1126 temperature simulations: Trends, long-term persistence, autocorrelation, and distributional shape, *Earth's*
1127 *Future*, 8, e2020EF001667, <https://doi.org/10.1029/2020EF001667>, 2020.

1128 Pascoe, C., Lawrence, B. N., Guilyardi, E., Jukes, M., and Taylor, K. E.: Documenting numerical experiments in
1129 support of the Coupled Model Intercomparison Project Phase 6 (CMIP6), *Geosci. Model Dev.*, 13, 2149–
1130 2167, <https://doi.org/10.5194/gmd-13-2149-2020>, 2020.

1131 Pendergrass, A. G., Gleckler, P. J., Leung, L. R., and Jakob, C.: Benchmarking simulated precipitation in earth system
1132 models, *B. Am. Meteorol. Soc.*, 101, E814–E816, <https://doi.org/10.1175/bams-d-19-0318.1>, 2020.

1133 Phillips, A. S., Deser, C., and Fasullo, J. T.: Evaluating modes of variability in climate models, *Eos, Transactions*
1134 *American Geophysical Union*, 95, 453–455, <https://doi.org/10.1002/2014eo490002>, 2014.

1135 Planton, Y., Guilyardi, E., Wittenberg, A. T., Lee, J., Gleckler, P. J., Bayr, T., McGregor, S., McPhaden, M. J., Power,
1136 S. B., Roehrig, R., Vialard, J., and Voltaire, A.: Evaluating Climate Models with the CLIVAR 2020 ENSO
1137 Metrics Package, *B. Am. Meteorol. Soc.*, 102, E193–E217, <https://doi.org/10.1175/bams-d-19-0337.1>, 2021.

Deleted: Bulletin of the American Meteorological Society

Deleted: Bulletin of the American Meteorological Society

Deleted: Geoscientific Model Development

Deleted: Bulletin of the American Meteorological Society

Deleted: Nature Communications

Deleted: Journal of Climate

Deleted: Bulletin of the American Meteorological Society

Deleted: Bulletin of the American Meteorological Society

- 1146 Planton, Y. Y., Lee, J., Wittenberg, A. T., Gleckler, P. J., Guilyardi, E., McGregor, S., and McPhaden, M. J.:
 1147 Estimating uncertainty in simulated ENSO statistics, *J. Adv. Model. Earth Sy.* (under review), ESS Open
 1148 Archive, <https://doi.org/10.22541/essoar.170196744.48068128/v1>, 2023. **Deleted: Journal of Advances in Modeling Earth Systems**
- 1149 Potter, G. L., Bader, D. C., Riches, M., Bamzai, A. and Joseph, R.: Celebrating two decades of the Program for Climate
 1150 Model Diagnosis and Intercomparison. *B. Am. Meteorol. Soc.*, 92, 629-631, **Deleted: Bulletin of the American Meteorological Society**
 1151 <https://doi.org/10.1175/2011BAMS3018.1>, 2011.
- 1152 Qin, Y., Zelinka, M. D., and Klein, S. A.: On the Correspondence Between Atmosphere-Only and Coupled
 1153 Simulations for Radiative Feedbacks and Forcing From CO₂, *J. Geophys. Res.-Atmos.*, 127, **Deleted: Journal of Geophysical Research: Atmospheres**
 1154 <https://doi.org/10.1029/2021jd035460>, 2022.
- 1155 Randall, D.A., Wood, R.A., Bony, S., Colman, R., Fichet, T., Fyfe, J., Kattsov, V., Pitman, A., Shukla, J., Srinivasan,
 1156 J. and Stouffer, R.J.: Climate models and their evaluation. In *Climate change 2007: The physical science
 1157 basis. Contribution of Working Group I to the Fourth Assessment Report of the IPCC (FAR)*, 589-662,
 1158 Cambridge University Press, 2007.
- 1159 Randerson, J. T., Hoffman, F. M., Thornton, P. E., Mahowald, N. M., Lindsay, K., Lee, Y.-H., Nevison, C. D., Doney,
 1160 S. C., Bonan, G. B., Stöckli, R., Covey, C., Running, S. W., and Fung, I.: Systematic assessment of terrestrial
 1161 biogeochemistry in coupled climate-carbon models, *Glob. Change Biol.*, 15, 2462–2484, **Deleted: Global Change Biology**
 1162 <https://doi.org/10.1111/j.1365-2486.2009.01912.x>, 2009.
- 1163 Rasch, P. J., Xie, S., Ma, P.-L., Lin, W., Wang, H., Tang, Q., Burrows, S. M., Caldwell, P., Zhang, K., Easter, R. C.,
 1164 Cameron-Smith, P. J., Singh, B., Wan, H., Golaz, J.-C., Harrop, B. E., Roesler, E. L., Bacmeister, J. T.,
 1165 Larson, V. E., Evans, K. J., Qian, Y., Taylor, M. A., Leung, L. R., Zhang, Y., Brent, L., Branstetter, M. L.,
 1166 Hannay, C., Mahajan, S., Mamejanov, A., Neale, R., Richter, J. H., Yoon, J.-H., Zender, C. S., Bader, D. C.,
 1167 Flanner, M., Foucar, J. G., Jacob, R., Keen, N. D., Klein, S. A., Liu, X., Salinger, A. G., Shrivastava, M., and
 1168 Yang, Y.: An overview of the atmospheric component of the Energy Exascale Earth System model, *J. Adv.
 1169 Model. Earth Sy.*, 11, 2377–2411, <https://doi.org/10.1029/2019ms001629>, 2019. **Deleted: Journal of Advances in Modeling Earth Systems**
- 1170 Reichler, T. and Kim, J.: How well do coupled models simulate today's climate?, *B. Am. Meteorol. Soc.*, 89, 303–
 1171 312, <https://doi.org/10.1175/bams-89-3-303>, 2008. **Deleted: Bulletin of the American Meteorological Society**
- 1172 Righi, M., Andela, B., Eyring, V., Lauer, A., Predoi, V., Schlund, M., Vegas-Regidor, J., Bock, L., Brötz, B., De
 1173 Mora, L., Diblen, F., Dreyer, L., Drost, N., Earnshaw, P., Hassler, B., Koldunov, N., Little, B., Tomas, S. L.,
 1174 and Zimmermann, K.: Earth System Model Evaluation Tool (ESMValTool) v2.0 – technical overview,
 1175 *Geosic. Model Dev.*, 13, 1179–1199, <https://doi.org/10.5194/gmd-13-1179-2020>, 2020. **Deleted: Geoscientific Model Development**
- 1176 Sanderson, B. M. and Wehner, M. F.: Weighting strategy for the Fourth National Climate Assessment, in: *Climate
 1177 Science Special Report: Fourth National Climate Assessment, Volume I*, edited by Wuebbles, D. J., Fahey,
 1178 D. W., Hibbard, K. A., Dokken, D. J., Stewart, B. C., and Maycock, T.K., U.S. Global Change Research
 1179 Program, Washington, DC, USA, 436-442, <https://doi.org/10.7930/J06T0JS3>, 2017.
- 1180 Sanderson, B. M., Wehner, M., and Knutti, R.: Skill and independence weighting for multi-model assessments,
 1181 *Geosic. Model Dev.*, 10, 2379–2395, <https://doi.org/10.5194/gmd-10-2379-2017>, 2017. **Deleted: Geoscientific Model Development**

1190 Sherwood, S. C., Webb, M. J., Annan, J. D., Armour, K. C., Forster, P. M., Hargreaves, J. C., Hegerl, G. C., Klein, S.
1191 A., Marvel, K., Rohling, E. J., Watanabe, M., Andrews, T., Braconnot, P., Bretherton, C. S., Foster, G. L.,
1192 Hausfather, Z., Von Der Heydt, A. S., Knutti, R., Mauritsen, T., Norris, J. R., Proistosescu, C., Rugenstein,
1193 M., Schmidt, G. A., Tokarska, K., and Zelinka, M. D.: An assessment of Earth's climate sensitivity using
1194 multiple lines of evidence, *Rev. Geophys.*, 58, <https://doi.org/10.1029/2019rg000678>, 2020.

1195 Singh, R., and AchutaRao, K.: Sensitivity of future climate change and uncertainty over India to performance-based
1196 model weighting. *Clim. Change*, <https://doi.org/10.1007/s10584-019-02643-y>, 2020.

1197 Sillmann, J., Kharin, V. V., Zhang, X., Zwiers, F. W., and Bronaugh, D.: Climate extremes indices in the CMIP5
1198 multimodel ensemble: Part I. Model evaluation in the present climate, *J. Geophys. Res.-Atmos.*, 118, 1716–
1199 1733, <https://doi.org/10.1002/jgrd.50203>, 2013.

1200 Sperber, K. R.: Madden-Julian variability in NCAR CAM2.0 and CCSM2.0, *Clim. Dyn.*, 23, 259–278,
1201 <https://doi.org/10.1007/s00382-004-0447-4>, 2004.

1202 Sperber, K. R., Annamalai, H., Kang, I.-S., Kitoh, A., Moise, A., Turner, A., Wang, B., and Zhou, T.: The Asian
1203 summer monsoon: an intercomparison of CMIP5 vs. CMIP3 simulation of the late 20th century. *Clim. Dyn.*,
1204 41, 2711–2744, <https://doi.org/10.1007/s00382-012-1607-6>, 2013.

1205 Sperber K. R., Gualdi, S., Legutke, S., Gayler, V.: The Madden–Julian oscillation in ECHAM4 coupled and uncoupled
1206 general circulation models, *Clim. Dyn.*, 25, 117–140, <https://doi.org/10.1007/s00382-005-0026-3>, 2005.

1207 Srivastava, A., Grotjahn, R., and Ullrich, P. A.: Evaluation of historical CMIP6 model simulations of extreme
1208 precipitation over contiguous US regions, *Weather and Climate Extremes*, 29, 100268,
1209 <https://doi.org/10.1016/j.wace.2020.100268>, 2020.

1210 Steed, C. A., Shipman, G., Thornton, P., Ricciuto, D., Erickson, D. and Branstetter, M.: Practical application of parallel
1211 coordinates for climate model analysis, *Procedia Comput. Sci.*, 9, 877–886,
1212 <https://doi.org/10.1016/j.procs.2012.04.094>, 2012.

1213 Stevens, B., Satoh, M., Auger, L., Biercamp, J., Bretherton, C. S., Chen, X., Düben, P., Judt, F., Khairoutdinov, M.,
1214 Klocke, D., Kodama, C., Kornbluh, L., Lin, S.-J., Neumann, P., Putman, W. M., Röber, N., Shibuya, R.,
1215 Vanniere, B., Vidale, P. L., Wedi, N., and Zhou, L.: DYAMOND: the DYnamics of the Atmospheric general
1216 circulation Modeled On Non-hydrostatic Domains, *Progress in Earth and Planetary Science*, 6,
1217 <https://doi.org/10.1186/s40645-019-0304-z>, 2019.

1218 Stoner, A. M. K., Hayhoe, K., and Wuebbles, D. J.: Assessing general circulation model simulations of atmospheric
1219 teleconnection patterns, *J. Climate*, 22, 4348–4372, <https://doi.org/10.1175/2009jcli2577.1>, 2009.

1220 Sung, H. M., Kim, J., Shim, S., Seo, J., Kwon, S.-H., Sun, M.-A., Moon, H.-J., Lee, J., Lim, Y. C., Boo, K.-O., Kim,
1221 Y., Lee, J., Lee, J., Kim, J.-S., Marzin, C., and Byun, Y.-H.: Climate change projection in the Twenty-First
1222 Century simulated by NIMS-KMA CMIP6 model based on new GHGs concentration pathways, *Asia-Pac. J.*
1223 *Atmos. Sci.*, 57, 851–862, <https://doi.org/10.1007/s13143-021-00225-6>, 2021.

1224 Tang, Q., Prather, M. J., Hsu, J., Ruiz, D. J., Cameron-Smith, P. J., Xie, S., and Golaz, J.-C.: Evaluation of the
1225 interactive stratospheric ozone (O3v2) module in the E3SM version 1 Earth system model, *Geosic. Model*
1226 *Dev.*, 14, 1219–1236, <https://doi.org/10.5194/gmd-14-1219-2021>, 2021.

Deleted: Reviews of Geophysics

Deleted: J

Formatted: Indent: Hanging: 0.5"

Deleted: Journal of Geophysical Research: Atmospheres

Deleted: Clim Dyn

Deleted: Clim Dyn

Deleted: Clim Dyn

Deleted: Procedia Computer Science

Deleted: Journal of Climate

Deleted: Asia-pacific Journal of Atmospheric Sciences

Deleted: Geoscientific Model Development

1237 Tang, S., Fast, J. D., Zhang, K., Hardin, J. C., Varble, A. C., Shilling, J. E., Mei, F., Zawadowicz, M. A., and Ma, P.-
1238 L.: Earth System Model Aerosol–Cloud Diagnostics (ESMAC Diags) package, version 1: assessing E3SM
1239 aerosol predictions using aircraft, ship, and surface measurements, *Geosci. Model Dev.*, 15, 4055–4076,
1240 <https://doi.org/10.5194/gmd-15-4055-2022>, 2022.

1241 Tang, S., Varble, A. C., Fast, J. D., Zhang, K., Wu, P., Dong, X., Mei, F., Pekour, M., Hardin, J. C., and Ma, P.-L.:
1242 Earth System Model Aerosol–Cloud Diagnostics (ESMAC Diags) package, version 2: assessing aerosols,
1243 clouds, and aerosol–cloud interactions via field campaign and long-term observations, *Geosci. Model Dev.*,
1244 16, 6355–6376, <https://doi.org/10.5194/gmd-16-6355-2023>, 2023.

1245 Taylor, K. E.: Truly Conserving with Conservative Remapping Methods, *Geosci. Model Dev.* 17, 415–430,
1246 <https://doi.org/10.5194/gmd-17-415-2024>, 2024.

1247 Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, *J. Geophys. Res.*, 106, 7183–
1248 7192, <https://doi.org/10.1029/2000jd900719>, 2001.

1249 Taylor, K. E., Stouffer, R. J., and Meehl, G. A.: An overview of CMIP5 and the experiment design, *B. Am. Meteorol.*
1250 *Soc.*, 93, 485–498, <https://doi.org/10.1175/bams-d-11-00094.1>, 2012.

1251 Teixeira, J., Waliser, D. E., Ferraro, R., Gleckler, P. J., Lee, T., and Potter, G. L.: Satellite observations for CMIP5:
1252 The Genesis of OBS4MIPs, *B. Am. Meteorol. Soc.*, 95, 1329–1334, <https://doi.org/10.1175/bams-d-12-00204.1>, 2014.

1253

1254 Tian, B., and Dong, X.: The Double-ITCZ Bias in CMIP3, CMIP5, and CMIP6 Models Based on Annual Mean
1255 Precipitation, *Geophys. Res. Lett.*, 47, e2020GL087232, <https://doi.org/10.1029/2020GL087232>, 2020

1256 Ullrich, P. A. and Zarzycki, C. M.: TempestExtremes: a framework for scale-insensitive pointwise feature tracking on
1257 unstructured grids, *Geosci. Model Dev.*, 10, 1069–1090, <https://doi.org/10.5194/gmd-10-1069-2017>, 2017.

1258 Ullrich, P. A., Zarzycki, C. M., McClenny, E., Pinheiro, M. C., Stansfield, A. M., and Reed, K. A.: TempestExtremes
1259 v2.1: a community framework for feature detection, tracking, and analysis in large datasets, *Geosci. Model*
1260 *Dev.*, 14, 5023–5048, <https://doi.org/10.5194/gmd-14-5023-2021>, 2021.

1261 U.S. Department of Energy (DOE): Benchmarking Simulated Precipitation in Earth System Models Workshop Report,
1262 DOE/SC-0203, U.S. Department of Energy Office of Science, Biological and Environmental Research (BER)
1263 Program. Germantown, Maryland, USA. 2020.

1264 Vo, T., Po-Chedley, P., Boutte, J., Zhang, C., Lee, J., Gleckler, P., Durack, P., Taylor, K., and Golaz, J.-C.: Xarray
1265 Climate Data Analysis Tools (xCDAT): A Python Package for Simple and Robust Analysis of Climate Data,
1266 The 103rd AMS Annual Meeting, Abstract, 2023.

1267 Waliser, D. E., Gleckler, P. J., Ferraro, R., Taylor, K. E., Ames, S., Biard, J., Bosilovich, M. G., Brown, O. B., Chepfer,
1268 H., Cinquini, L., Durack, P. J., Eyring, V., Mathieu, P.-P., Lee, T., Pinnock, S., Potter, G. L., Rixen, M.,
1269 Saunders, R., Schulz, J. B., Thépaut, J.-N., and Tuma, M.: Observations for Model Intercomparison Project
1270 (Obs4MIPs): status for CMIP6, *Geosci. Model Dev.*, 13, 2945–2958, [https://doi.org/10.5194/gmd-13-2945-](https://doi.org/10.5194/gmd-13-2945-2020)
1271 2020, 2020.

1272 Waliser, D. E., Sperber, K. R., Hendon, H. H., Kim, D., Maloney, E. D., Wheeler, M. C., Weickmann, K. M., Zhang,
1273 C., Donner, L. J., Gottschalck, J., Higgins, W., Kang, I. S., Legler, D. M., Moncrieff, M. W., Schubert, S. D.,

Deleted: Discuss. [preprint]
Deleted: gmd-2023-177
Deleted: in review,
Deleted: 3
Deleted: Journal of Geophysical Research
Deleted: Bulletin of the American Meteorological Society

Deleted: Bulletin of the American Meteorological Society

Deleted: Geophysical Research Letters

Deleted: Geoscientific Model Development

Deleted: Geoscientific Model Development

Deleted: Geoscientific Model Development

1285 Stern, W., Vitart, F., Wang, B., Wang, W., and Woolnough, S. J.: MJO Simulation Diagnostics, *J. Climate*,
1286 22, 3006–3030, <https://doi.org/10.1175/2008jcli2731.1>, 2009.

1287 Wang, J., Liu, X., Shen, H. W. and Lin, G.: Multi-resolution climate ensemble parameter analysis with nested parallel
1288 coordinates plots, *IEEE T. Vis. Comput. G. R.*, 23, 81-90, <https://doi.org/10.1109/TVCG.2016.2598830>,
1289 2017.

1290 Wehner, M., Gleckler, P. J., and Lee, J.: Characterization of long period return values of extreme daily temperature
1291 and precipitation in the CMIP6 models: Part 1, model evaluation, *Weather and Climate Extremes*, 30,
1292 100283, <https://doi.org/10.1016/j.wace.2020.100283>, 2020.

1293 Wehner, M., Lee, J., Risser, M. D., Ullrich, P. A., Gleckler, P. J., and Collins, W. D.: Evaluation of extreme sub-daily
1294 precipitation in high-resolution global climate model simulations, *Philos. T. R. Soc. A.*, 379, 20190545,
1295 <https://doi.org/10.1098/rsta.2019.0545>, 2021.

1296 Whitehall, K., Mattmann, C., Waliser, D., Kim, J., Goodale, C., Hart, A., Ramirez, P., Zimdars, P., Crichton, D.,
1297 Jenkins, G., Jones, C., Asrar, G., and Hewitson, B.: Building Model Evaluation and Decision Support
1298 Capacity for CORDEX, *WMO Bulletin*, 61, available at:
1299 [https://public.wmo.int/en/resources/bulletin/building-model-evaluation-and-decision-support-capacity-](https://public.wmo.int/en/resources/bulletin/building-model-evaluation-and-decision-support-capacity-cordex)
1300 [cordex](https://public.wmo.int/en/resources/bulletin/building-model-evaluation-and-decision-support-capacity-cordex) (last access date: 14 September 2023), 2012.

1301 Williams, D. N.: Visualization and analysis tools for ultrascale climate data, *Eos, Transactions American Geophysical*
1302 *Union*, 95, 377–378, <https://doi.org/10.1002/2014eo420002>, 2014.

1303 Williams, D. N., Balaji, V., Cinquini, L., Denvil, S., Duffy, D. Q., Evans, B., Ferraro, R., Hansen, R., Lautenschlager,
1304 M., and Trenham, C.: A global repository for Planet-Sized experiments and observations, *B. Am. Meteorol.*
1305 *Soc.*, 97, 803–816, <https://doi.org/10.1175/bams-d-15-00132.1>, 2016.

1306 Williams, D. N., Doutriaux, C., Drach, R., and McCoy, R.: The Flexible Climate Data Analysis Tools (CDAT) for
1307 Multi-model Climate Simulation Data, *IEEE International Conference on Data Mining Workshops*, 254–261,
1308 <https://doi.org/10.1109/icdmw.2009.64>, 2009.

1309 Wong, P. C., Shen, H. W., Leung, R., Hagos, S., Lee, T. Y., Tong, X. and Lu, K.: Visual analytics of large-scale
1310 climate model data, in: 2014 IEEE 4th Symposium on Large Data Analysis and Visualization (LDAV), 85-
1311 92, <https://doi.org/10.1109/LDAV.2014.7013208>, 2014.

1312 Xie, P., Joyce, R., Wu, S., Yoo, S.H., Yarosh, Y., Sun, F. and Lin, R.: Reprocessed, bias-corrected CMORPH global
1313 high-resolution precipitation estimates from 1998, *J. Hydrometeorol.*, 18, 1617-1641, 2017.

1314 Xue, Z. and Ullrich, P. A.: A Comprehensive Intermediate-Term Drought Evaluation System and Evaluation of
1315 Climate Data Products over the Conterminous United States, *J. Hydrometeorol.*, [https://doi.org/10.1175/jhm-](https://doi.org/10.1175/jhm-d-20-0314.1)
1316 [d-20-0314.1](https://doi.org/10.1175/jhm-d-20-0314.1), 2021.

1317 Young, A. H., Knapp, K. R., Inamdar, A. K., Hankins, W., and Rossow, W. B.: The International Satellite Cloud
1318 Climatology Project H-Series climate data record product, *Earth Syst. Sci. Data*, 10, 583–593,
1319 <https://doi.org/10.5194/essd-10-583-2018>, 2018.

1320 Zelinka, M. D., Klein, S. A., Qin, Y., and Myers, T. A.: Evaluating climate models’ cloud feedbacks against expert
1321 judgment, *J. Geophys. Res.-Atmos.*, 127, <https://doi.org/10.1029/2021jd035198>, 2022.

Deleted: Journal of Climate

Deleted: IEEE Transactions on Visualization and Computer Graphics...

Deleted: Philosophical Transactions of the Royal Society A...

Deleted: Bulletin of the American Meteorological Society

Deleted: Journal of Hydrometeorology

Deleted: Journal of Hydrometeorology

Deleted: Earth System Science Data

Deleted: Journal of Geophysical Research: Atmospheres

1332 Zelinka, M. D., Myers, T. A., McCoy, D. T., Po-Chedley, S., Caldwell, P. M., Ceppi, P., Klein, S. A., and Taylor, K.
1333 E.: Causes of higher climate sensitivity in CMIP6 models, [Geophys. Res. Lett.](https://doi.org/10.1029/2019GL085782), 47, e2019GL085782,
1334 <https://doi.org/10.1029/2019GL085782>, 2020.

1335 Zhang, C., Golaz, J.-C., Forsyth, R., Vo, T., Xie, S., Shaheen, Z., Potter, G. L., Asay-Davis, X. S., Zender, C. S., Lin,
1336 W., Chen, C.-C., Terai, C. R., Mahajan, S., Zhou, T., Balaguru, K., Tang, Q., Tao, C., Zhang, Y.,
1337 Emmenegger, T., Burrows, S., and Ullrich, P. A.: The E3SM Diagnostics Package (E3SM Diags v2.7): a
1338 Python-based diagnostics package for Earth system model evaluation, *Geosci. Model Dev.*, 15, 9031–9056,
1339 <https://doi.org/10.5194/gmd-15-9031-2022>, 2022.

1340 Zhang, C. and Hendon, H. H.: Propagating and standing components of the intraseasonal oscillation in tropical
1341 convection, *J. Atmos. Sci.*, 54, 741–752, [https://doi.org/10.1175/1520-0469\(1997\)054](https://doi.org/10.1175/1520-0469(1997)054), 1997.

1342 Zhang, C., Xie, S., Klein, S. A., Ma, H. Y., Tang, S., Van Weverberg, K., Morcrette, C. J. and Petch, J.: CAUSES:
1343 Diagnosis of the summertime warm bias in CMIP5 climate models at the ARM Southern Great Plains site.
1344 *J. Geophys. Res.-Atmos.*, 123, 2968-2992, <https://doi.org/10.1002/2017JD027200>, 2018.

1345 Zhang, C., Xie, S., Tao, C., Tang, S., Emmenegger, T., Neelin, J. D., Schiro, K. A., Lin, W. and Shaheen, Z.: The
1346 ARM data-oriented metrics and diagnostics package for climate models: A new tool for evaluating climate
1347 models with field data, <https://doi.org/10.1175/BAMS-D-19-0282.1>, *B. Am. Meteorol. Soc.*, 101, E1619-
1348 E1627, 2020.

1349 Zhao, B., Lin, P., Hu, A., Liu, H., Ding, M., Yu, Z., and Yu, Y.: Uncertainty in Atlantic Multidecadal Oscillation
1350 derived from different observed datasets and their possible causes, *Frontiers in Marine Science*, 9,
1351 <https://doi.org/10.3389/fmars.2022.1007646>, 2022.

1352 Zhao, M., Golaz, J.-C., Held, I. M., Guo, H., Balaji, V., Benson, R., Chen, J. H., Chen, X., Donner, L. J., Dunne, J.,
1353 Dunne, K. A., Durachta, J., Fan, S.-M., Freidenreich, S. M., Garner, S. T., Ginoux, P., Harris, L., Horowitz,
1354 L. W., Krasting, J. P., Langenhorst, A. R., Zhi, L., Lin, P., Lin, S. J., Malyshev, S., Mason, E., Milly, P. C.
1355 D., Ming, Y., Naik, V., Paulot, F., Paynter, D., Philipps, P. J., Radhakrishnan, A., Ramaswamy, V.,
1356 Robinson, T., Schwarzkopf, D., Seman, C. J., Shevliakova, E., Shen, Z., Shin, H. H., Silvers, L. G., Wilson,
1357 J. R., Winton, M., Wittenberg, A. T., Wyman, B., and Xiang, B.: The GFDL Global Atmosphere and Land
1358 Model AM4.0/LM4.0: 1. Simulation characteristics with prescribed SSTs, *J. Adv. Model. Earth Sy.*, 10,
1359 691–734, <https://doi.org/10.1002/2017ms001208>, 2018.

1360

Deleted: Geophysical Research Letters

Deleted: Journal of the Atmospheric Sciences

Deleted: Journal of Geophysical Research: Atmospheres

Deleted: Bulletin of the American Meteorological Society

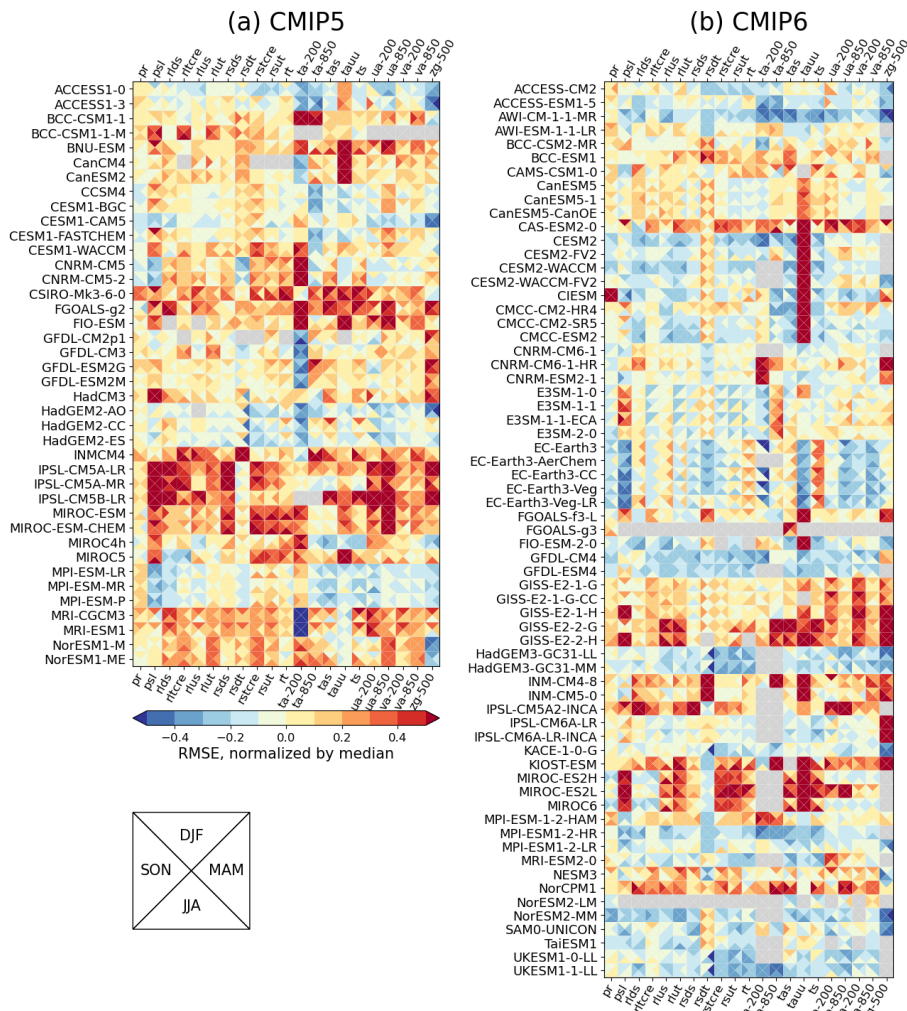
Deleted: Journal of Advances in Modeling Earth Systems

1366
1367
1368
1369

Table 1. List of variables and observation datasets used as reference datasets for the PMP's mean climate evaluation in this paper (Sect. 3.1 and Figs. 1-2). A ditto mark (") indicates the same as above.

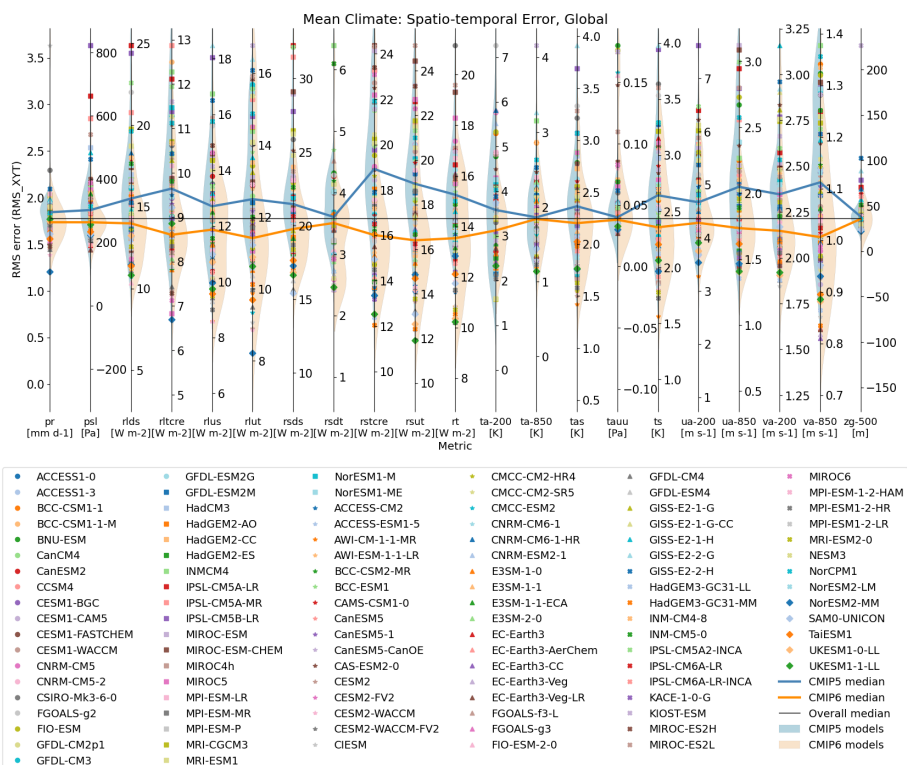
Deleted: Section

Variable	Variable full name	Product	Reference
ps	Precipitation	GPCP-2-3	Adler et al. (2018)
psl	Sea level pressure	ERA-5	Hersbach et al. (2020)
rlds	Surface Downwelling Longwave Radiation	CERES-EBAF-4-1	Loeb et al. (2018)
rltcre	Longwave cloud radiative effect	"	
rlus	Surface Upwelling Longwave Radiation	"	
rlut	Upwelling longwave at the top of atmosphere	"	
rsds	Surface Downwelling Shortwave Radiation	"	
rsdt	TOA Incident Shortwave Radiation	"	
rstcre	Shortwave cloud radiative effect	"	
rsut	Upwelling shortwave at the top of atmosphere	"	
rt	Net radiative flux	"	
ta-200, ta-850	Air temperature at 850 and 200 hPa	ERA-5	Hersbach et al. (2020)
tas	2-m air temperature	"	
tauu	Surface zonal wind stress	ERA-INT	Dee et al. (2011)
ts	Surface temperature	ERA-5	Hersbach et al. (2020)
ua-200, ua-850	Zonal wind component at 850 and 200 hPa	"	
va-200, va-850	Meridional wind component at 850 and 200 hPa	"	
zg-500	Geopotential height at 500 hPa	"	



1371
 1372 **Figure 1.** Portrait plot for spatial RMSE (uncentered) of global seasonal climatologies for (a)
 1373 CMIP5 (models ACCESS1-0 to NorESM1-ME on the ordinate) and (b) CMIP6 (models
 1374 ACCESS-CM2 to UKESM1-1-LL on the ordinate) for 1981-2005 epoch. The RMSE is calculated
 1375 for each season (shown as triangles in each box) over the globe including both land and ocean,
 1376 and model and reference data were interpolated to a common 2.5x2.5 degree grid. The RMSE
 1377 of each variable is normalized by the median RMSE of all CMIP5 and 6 models. A result of 0.2
 1378 (-0.2) is indicative of an error that is 20% greater (lesser) than the median RMSE across all
 1379 models. Models in each group are sorted in alphabetical order. Full names of variable names on
 1380 the abscissa and their reference datasets can be found in Table 1. Detailed information for

1381 models can be found at the *Earth System Documentation* (ES-DOC, <https://search.es-doc.org/>;
1382 Pascoe et al., 2020). The interactive version of the Portrait plot in this figure is available on the
1383 PMP result pages on the PCMDI website (https://pcmdi.llnl.gov/metrics/mean_clim/).



1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396

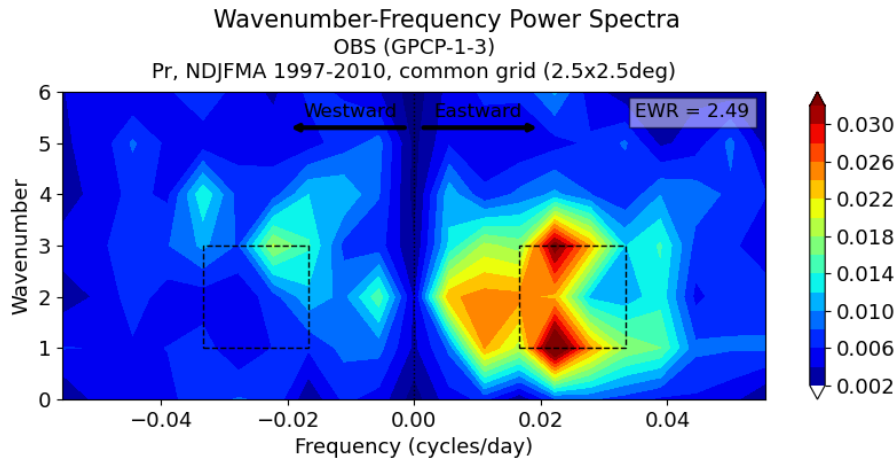
Figure 2. Parallel Coordinate Plot for spatio-temporal RMSE (Gleckler et al., 2008) from mean climate evaluation. Each vertical axis represents a different variable. Results from each model are displayed as symbols. Middle of each vertical axis is aligned with the median statistic of all CMIP5 and CMIP6 models. The cross-generation model distributions of model performance are shaded on the left (CMIP5, blue) and right (CMIP6, orange) sides of each axis. Also, medians from CMIP5 (blue) and CMIP6 (orange) model groups are highlighted as lines. Full names for model variables on the abscissa and their reference datasets can be found in Table 1. Time epoch used for this analysis is 1981-2005. Detailed information for models can be found at the *Earth System Documentation* (ES-DOC, <https://search.es-doc.org/>; Pascoe et al., 2020). The interactive version of the Portrait plot in this figure is available on the PMP result pages on the PCMDI website (https://pcmdi.llnl.gov/metrics/mean_clim/).



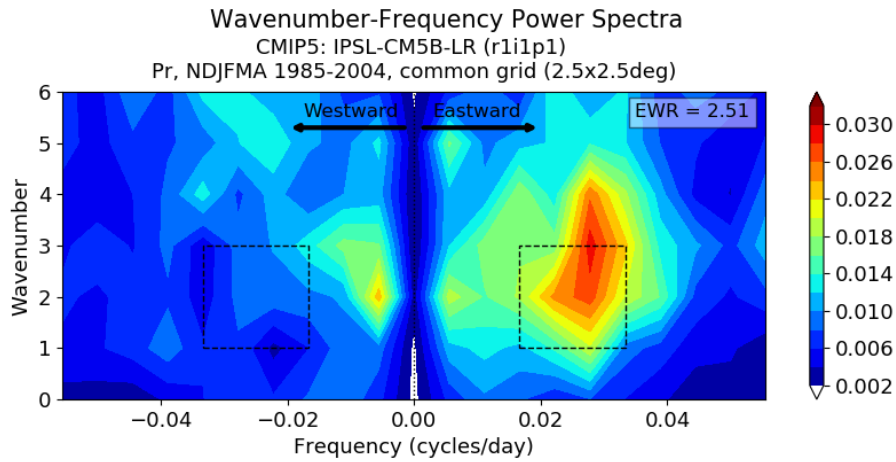
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407

Figure 3. Application of ENSO metrics to CMIP6 models. Model names with an asterisk (*) indicate that 10 or more ensemble members were used in this analysis. Dots indicate metric values from individual ensemble members while bars indicate the average of metric values across the ensemble members. Bars colored for easier identification of model names at the bottom of the figure. Metrics were grouped into three *Metric Collections*: (a-n) ENSO Performance, (o-r) ENSO Teleconnections, and (s-w) ENSO processes. Names of individual metrics and default reference datasets being used are noted on top of each panel, and observational uncertainty by applying the metrics for alternative reference datasets noted on the upper right of each panel is shown as gray-shaded. Detailed descriptions for each metric can be found at https://github.com/CLIVAR-PRP/ENSO_metrics/wiki.

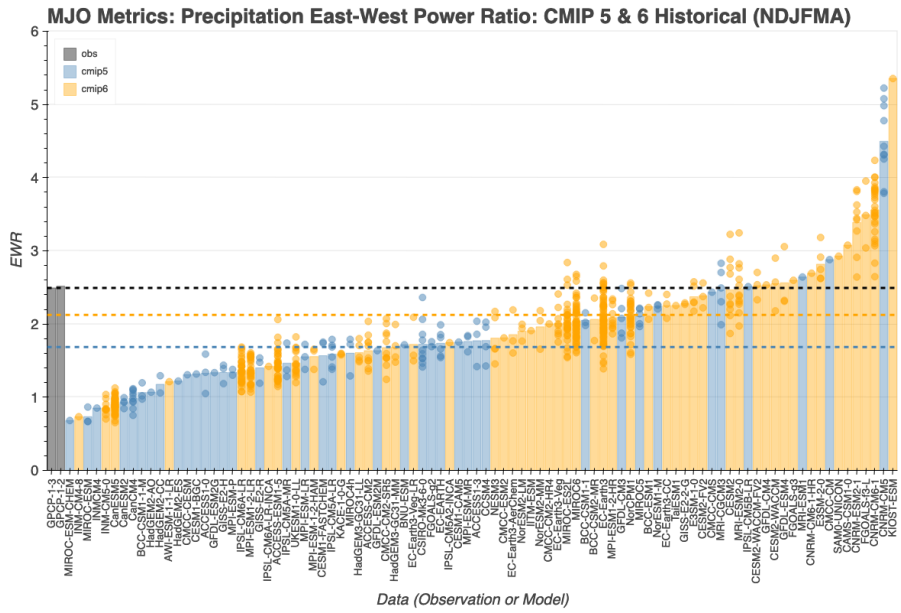
1434
1435 (a) Observation



1436
1437 (b) Model

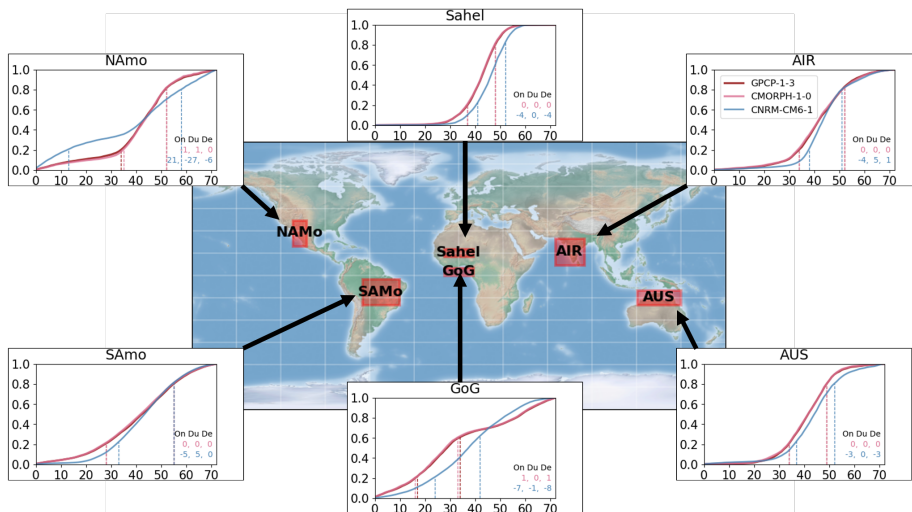


1438
1439
1440 **Figure 5.** MJO EWR diagnostics – wavenumber-frequency power spectra – from (a) GPCP v1.3
1441 (Huffman et al., 2001) and (b) IPSL-CM5B-LR model of CMIP5. The EWR is defined as the ratio
1442 of eastward power (averaged in the box on the right) to westward power (averaged in the box
1443 on the left) from the 2-dimensional wavenumber-frequency power spectra of daily 10°S–10°N
1444 averaged precipitation in November to April (shaded, $\text{mm}^2 \text{day}^{-2}$). Power spectra are calculated
1445 for each year and then averaged over all years of data. The units of power spectra for the
1446 precipitation is $\text{mm}^2 \text{day}^{-2}$ per frequency interval per wavenumber interval.



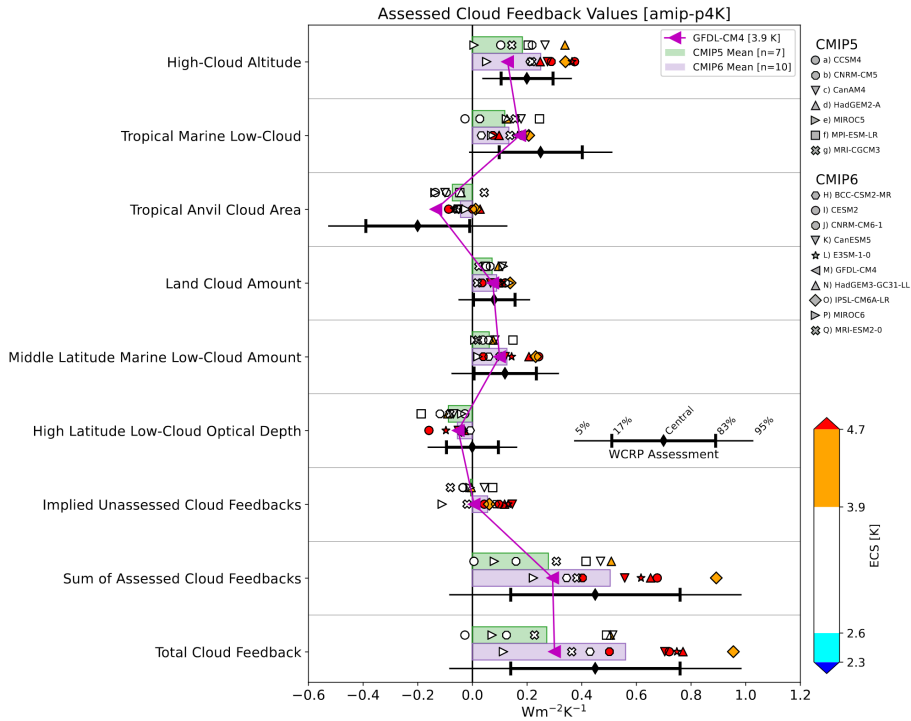
1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

Figure 6. MJO East-West Power Ratio (EWR, *unitless*) from CMIP5 and CMIP6 models, models in two different groups (CMIP5: blue, CMIP6: orange) are sorted by the value of the metric and compared to two observation datasets (purple, GPCP v1.2 and v1.3; Huffman et al., 2001). Horizontal dashed lines indicate EWR from the default primary reference observation (i.e., GPCP v1.3, black), averages of CMIP5 and CMIP6 models. The interactive plot is available at <https://pcmdi.llnl.gov/research/metrics/mjo/> where the horizontal axis can be resorted by CMIP group or model names as well. Hover mouse over boxes will show tooltips for metric values and a preview of dive-down plots that are shown in Figure 5.



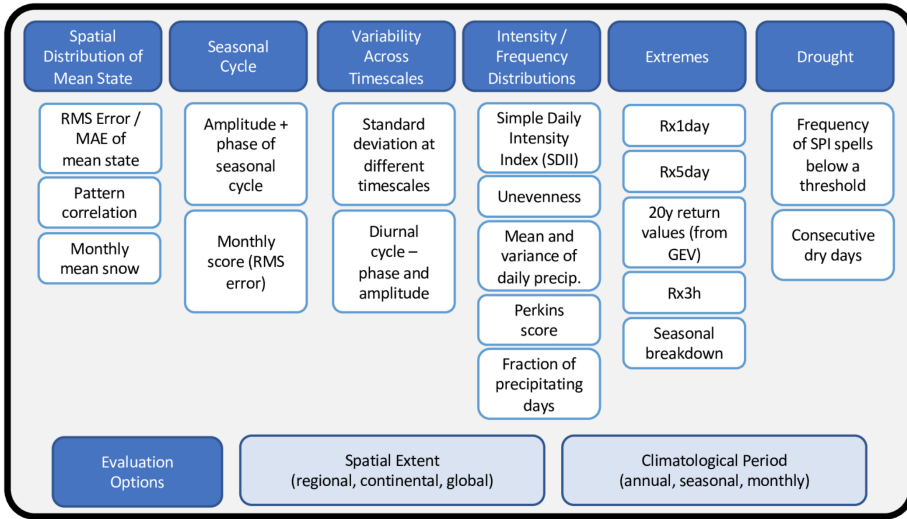
1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467

Figure 7. Demonstration of the monsoon metrics obtained from observation datasets (GPCP v1.3 and CMORPH v1.0 (Joyce et al., 2004; Xie et al., 2017)) and a CMIP6 model's Historical simulation conducted using CNRM-CM6-1. The results are obtained for monsoon regions: All-India Rainfall (AIR), Sahel, Gulf of Guinea (GoG), North American Monsoon (NAM), South American Monsoon (SAM), and Northern Australia (AUS). The regions are defined in Sperber and Annamalai (2014). Metrics for onset (On), Duration (Du), and Decay (De) derived as differences to the default observation (GPCP v1.3) in pentad indices (observation minus model) are shown at lower right of each panel. Pentad indices for onset and decay of each region are also shown as vertical lines.



1469 **Figure 8.** Cloud feedback components estimated in amp-p4K simulations from CMIP5 and
 1470 CMIP6 models. Symbols indicate individual model values, while horizontal bars indicate multi-
 1471 model means. Each model is color-coded by its ECS, with color boundaries corresponding to the
 1472 likely and very likely ranges of ECS as determined in Sherwood et al (2020). Each component's
 1473 expert-assessed likely and very likely confidence intervals are indicated with black error bars. An
 1474 illustrative model (GFDL-CM4) is highlighted.

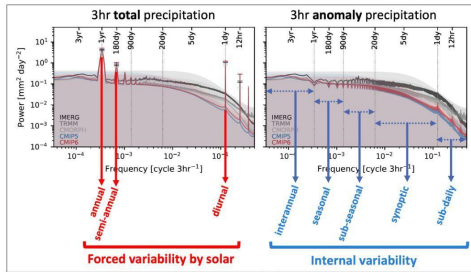
1476



1477
1478
1479
1480
1481
1482
1483

Figure 9. Proposed suite of baseline metrics for simulated precipitation benchmarking (figure reprinted from workshop report; US DOE, 2020).

(a) Power spectra (Tropics)

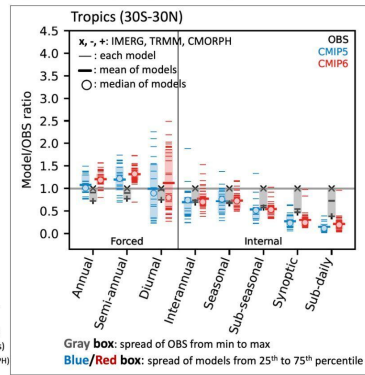


$$\text{Metric} = P_{\text{MODEL}}/P_{\text{OBS}}$$

P: selected or band-averaged power

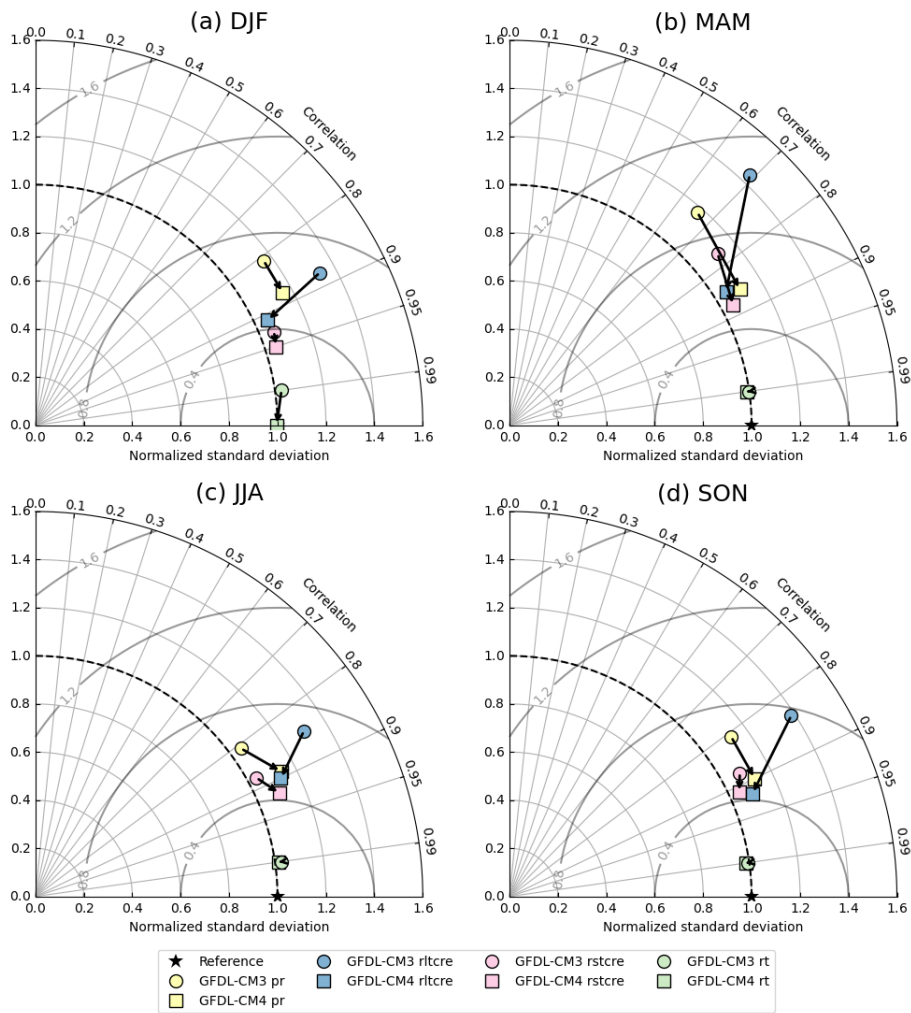
21 CMIP5 (53 realizations)
 33 CMIP6 (143 realizations)
 3 OBS (IMERG, TRMM, CMORPH)

(b) Metric for precip variability across timescales



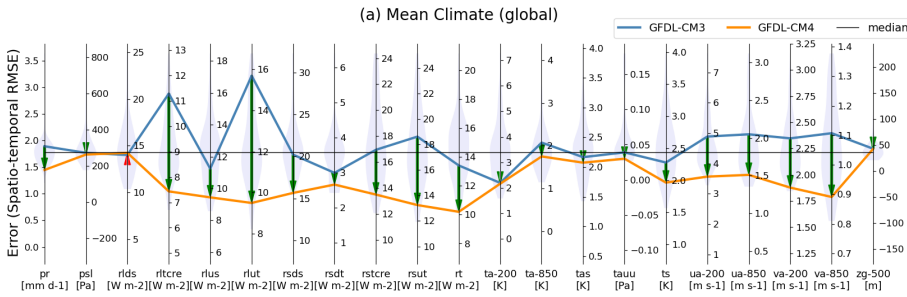
1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498

Figure 10. Example (a) an underlying diagnostic and (b) its resulting metrics for precipitation variability across timescales. (a) Power spectra of 3-hourly total (left) and anomaly (right) precipitation from IMERG (black), TRMM (gray), CMORPH (silver), CMIP5 (blue), and CMIP6 (red) averaged over the tropics (30°S-30°N). The colored shading indicates the 95% confidence interval for each observational product and for the CMIP5 and CMIP6 means. (b) Metrics for forced and internal precipitation variability based on power spectra. The reference observational product displayed is GPM IMERG (Huffman et al., 2015). The gray boxes represent the spread of the three observational products (“X” for IMERG, “-” for TRMM, and “+” for CMORPH) from the minimum to maximum values. Blue and red boxes indicate the 25th to the 75th percentile of CMIP models as a measure of spread. Individual models are shown as thin dashes, the multimodel mean as a thick dash, and the multimodel median as an open circle. Details for the diagnostics and metrics are described in Ahn et al. (2022).

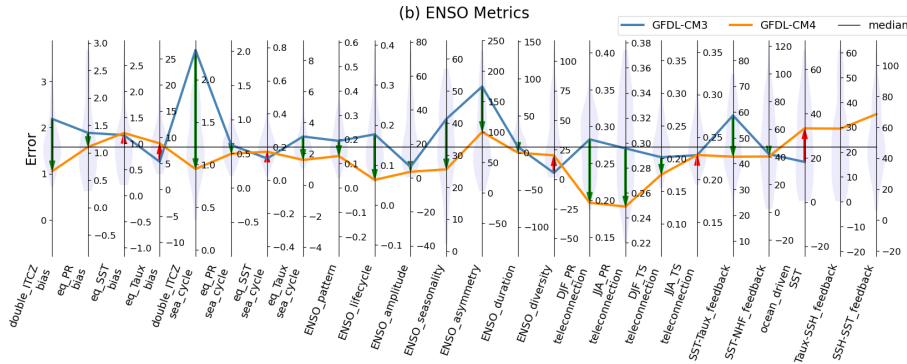


1499
 1500 **Figure 11.** Taylor Diagram contrasting performance of an ESM in their two different versions
 1501 (i.e., GFDL-CM3 from CMIP5 and GFDL-CM4 from CMIP6) in their Historical simulation for
 1502 multiple variables (precipitation [pr], longwave cloud radiative effect [rltcre], shortwave cloud
 1503 radiative effect [rstcre], and total radiation flux [rt]) in their climatology over the globe for (a) DJF,
 1504 (b) MAM, (c) JJA and (d) SON seasons. The arrow is directed toward the newer version of the
 1505 model from the older version (i.e., GFDL-CM3 → GFDL-CM4).

1506



1507



1508

1509

1510

1511

1512

1513

1514

1515

1516

1517

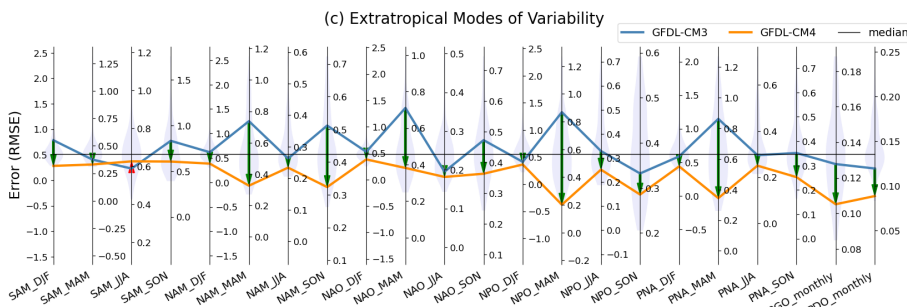


Figure 12. Parallel Coordinate Plot contrasting performance of two different versions of the GFDL model (i.e., GFDL-CM3 from CMIP5 and GFDL-CM4 from CMIP6) in their Historical experiment for errors from (a) mean climate, (b) ENSO, and (c) extratropical modes of variability. Improvement (degradation) in the later version of the model is highlighted as a downward green (upward red) arrow between lines. Middle of each vertical axis is set to the median of the group of benchmarking models (i.e., CMIP5 and CMIP6), with the axis range stretched to maximum distance to either minimum or maximum from the median for visual consistency. The inter-model distributions of model performance are shown as shaded violin plots along each vertical axis.