

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 [Geoscientific Model Development \(GMD\)](#) 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

31 **Abstract**

32

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

50 **1 Introduction**

51 Earth System Models (ESMs) are key tools for projecting climate change and conducting research to enhance
52 our understanding of the Earth system. With the advancements in computing power and the increasing importance of
53 climate projections, there has been an exponential growth of data size and diversity of ESM simulations. During the
54 1990's, the Atmospheric Model Intercomparison Project (AMIP; Gates, 1992; Gates et al., 1999) was a centralizing
55 activity within the modeling community, which led to the creation of the Coupled Model Intercomparison Project
56 (CMIP; Meehl et al., 1997, 2000, 2007; Covey et al., 2003; Taylor et al., 2012). Since 1989, the Program for Climate
57 Model Diagnosis and Intercomparison (PCMDI) has worked closely with the World Climate Research Programme's
58 (WCRP) Working Group on Coupled Models (WGCM) and Working Group on Numerical Experimentation (WGNE)
59 to design and implement these projects (Potter et al., 2011). The most recent phase of CMIP (CMIP6; Eyring et al.,
60 2016) provides a set of well-defined experiments that most climate modeling centers perform, and subsequently makes
61 results available for a large and diverse community to analyze.

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

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

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

95 In parallel, the current progress towards systematic model evaluation remains dynamic, with evolving
96 approaches and many independent paths being pursued. This has resulted in the development of diversified model
97 evaluation software packages. Examples in addition to the PMP include the ESMValTool (Eyring et al., 2016, 2019,
98 2020; Righi et al., 2020), the Model Diagnostics Task Force (MDTF) Diagnostics package (Maloney et al., 2019;
99 Neelin et al., 2023), the International Land Model Benchmarking (ILAMB) Software System (Collier et al., 2018) that
100 focuses on land surface and carbon cycle metrics, and the International Ocean Model Benchmarking (IOMB) Software
101 System (Fu et al., 2022) that focuses on surface and upper ocean biogeochemical variables. Some tools have been
102 developed with a more targeted focus on a specific subject area, such as the Climate Variability Diagnostics Package
103 (CVDP) that diagnoses climate variability modes (Phillips et al., 2014; Fasullo et al., 2020), and the Analyzing Scales
104 of Precipitation (ASoP) that focuses on analyzing precipitation scales across space and time (Klingaman et al., 2017;
105 Martin et al., 2017; Ordonez et al., 2021). The regional climate community also has actively developed metrics
106 packages such as the Regional Climate Model Evaluation System (RCMES; Lee et al., 2018a; Whitehall et al. 2012).
107 Separately, a few climate modeling centers have developed their own model evaluation packages to assist in their in-
108 house ESM development, e.g., the E3SM Diags (Zhang et al., 2022). There also have been other efforts to enhance
109 the usability of in-situ and field campaign observations in ESM evaluations, such as Atmospheric Radiation
110 Measurement (ARM) GCM Diag (Zhang et al., 2018, 2020) and Earth System Model Aerosol–Cloud Diagnostics
111 (ESMAC Diags; Tang et al., 2022, 2023). While they all have their own scientific priorities and technical approaches,
112 the uniqueness of the PMP is its focus on the objective characterization of the physical climate system as simulated
113 by community models. An important prioritization of the PMP is to advance all aspects of its workflow, in an open,
114 transparent, and reproducible manner, which is critical for benchmarking. The PMP summary statistics characterizing
115 CMIP simulations are version-controlled and made publicly available as a resource to the community.

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

124

125 **2 Software package and data description**

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

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

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

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

161 For reference datasets, the PMP uses observational products processed to be compliant with the Observations
162 for Model Intercomparison Projects (obs4MIPs; <https://pcmdi.github.io/obs4MIPs/>). The obs4MIPs effort was
163 initiated circa 2010 (Gleckler et al., 2011) to advance the use of the observations in model evaluation and research.
164 Substantial progress has been made in establishing obs4MIPs data standards that technically align with CMIP model
165 output (e.g., Teixeira et al., 2014; Ferraro et al., 2015), with the data products published on the ESGF (Waliser et al.,
166 2020). Obs4MIPs-compliant data were prepared with CMOR, and the data directly available via obs4MIPs are used
167 as PMP reference datasets.

168 The PMP leverages other Python-based open-source tools and libraries such as *xarray* (Hoyer and Hamman,
169 2017), *eofs* (Dawson, 2016), and many others. One of the primary fundamental tools used in the latest PMP version
170 is the Python package, Xarray Climate Data Analysis Tools (xCDAT; Vo et al., 2023; <https://xcdat.readthedocs.io>).
171 The xCDAT is developed to provide a more efficient, robust, and streamlined user experience in climate data analysis
172 when using *xarray* (<https://docs.xarray.dev/>). Portions of the PMP rely on the precursor of the xCDAT, a Python
173 library called Community Data Analysis Tools (CDAT, Williams et al., 2009; Williams, 2014; Doutriaux et al., 2019),
174 which has been fundamental since the early development stages of the PMP. The *xarray* software provides much of
175 the functionality of CDAT (e.g., I/O, indexing, and subsetting). However, it lacks some key climate domain features
176 that have been frequently used by scientists and exploited by the PMP (e.g., regridding, utilization of spatial/temporal
177 bounds for computational operations) which motivated the development of the xCDAT. Completing the transition
178 from CDAT to xCDAT is a technical priority for the next version of PMP.

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

196

197 **3 Current PMP capabilities**

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

207 208 **3.1 Climatology**

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

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

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

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

268

269 **3.2 El Niño-Southern Oscillation**

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

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

299
300 **3.3 Extratropical Modes of Variability**

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

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

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

324 The PMP's ETMoV metrics have been used in several model evaluation studies. For example, Orbe et al.
325 (2020) analyzed models from U.S. climate modeling groups including the U.S. Department of Energy (DOE), National
326 Aeronautics and Space Administration (NASA), National Center for Atmospheric Research (NCAR), and National
327 Oceanic and Atmospheric Administration (NOAA), where they found that the improvement in the ETMoV
328 performance is highly dependent on mode and season, when comparing across different generations of those models.
329 Sung et al. (2021) examined the performance of models run at the Korea Meteorological Administration (K-ACE and
330 UKESM1) in reproducing ETMoVs from their Historical simulations, and concluded that these models reasonably
331 capture most ETMoVs. Lee et al. (2021b) collectively evaluated ~130 models from CMIP3, 5, and 6 archive databases
332 using their ~850 Historical and ~300 AMIP simulations, where they found the spatial pattern skill improved in CMIP6
333 compared to CMIP5 or CMIP3 for most modes and seasons, while the improvement in amplitude skill is not clear.
334 Arcodia et al. (2023) used the PMP to derive PDO and AMO to investigate their role in decadal variability of
335 subseasonal predictability of precipitation over the western coast of North America and concluded that no significant
336 relationship was found.

337

338 ***3.4 Intraseasonal Oscillation***

339 The PMP has implemented metrics for the Madden-Julian Oscillation (MJO; Madden and Julian, 1971, 1972,
340 1994). The MJO is the dominant mode of tropical intraseasonal variability, characterized by a pronounced eastward
341 propagation of large-scale atmospheric circulation coupled with convection with a typical periodicity of 30-60 days.

342 Selected metrics from the MJO diagnostics package, developed by the CLIVAR MJO Working Group (Waliser et al.,
343 2009), have been implemented in the PMP following Ahn et al. (2017).

344 We have particularly focused on metrics for the MJO propagation: East/West power Ratio (EWR) and East
345 power normalized by Observation (EOR). The EWR is proposed by Zhang and Hendon (1997) which is defined as
346 the ratio of the total spectral power over the MJO band (eastward propagating, wavenumber 1-3 and period of 30-60
347 days) to that of its westward propagating counterpart in the wavenumber-frequency power spectra. The EWR metric
348 has been widely used in the community, to examine the robustness of the eastward propagating feature of the MJO
349 (e.g., Hendon et al., 1999; Lin et al., 2006; Kim et al., 2009; Ahn et al., 2017). The EOR is formulated by normalizing
350 a model's spectral power within the MJO band by the corresponding observed value. Ahn et al. (2017) showed EWRs
351 and EORs of the CMIP5 models. Using daily precipitation, the PMP calculates EWR and EOR separately for boreal
352 winter (November to April) and boreal summer (March to October). We apply the frequency-wavenumber
353 decomposition method to precipitation from observations (GPCP-based; 1997-2010) and the CMIP5 and CMIP6
354 Historical simulations for 1985-2004. For disturbances with wavenumbers 1-3 and frequencies corresponding to 30-
355 60 days, it is clear in observations that the eastward propagating signal dominates over its westward propagating
356 counterpart with an EWR value of about 2.49 (Fig. 5a). Figure 5b shows the wavenumber-frequency power spectrum
357 from CMIP5 IPSL-CM5B-LR as an example, which has an EWR value that is comparable to the observed value.

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

371 372 **3.5 Monsoons**

373 Based on the work of Sperber and Annamalai (2014), skill metrics in the PMP quantify how well models
374 represent the onset, decay, and duration of regional monsoons. From observations and Historical simulations, the
375 climatological pentads of precipitation are area-averaged for six monsoon-related domains: All-India Rainfall, Sahel,
376 Gulf of Guinea, North American Monsoon, South American Monsoon, and Northern Australia, as seen in Fig. 7. For
377 the domains in the Northern Hemisphere, the 73 climatological pentads run from January to December, while for the
378 domains in the Southern Hemisphere, the pentads run from July to June. For each domain, the precipitation is

379 accumulated at each subsequent pentad and then divided by the total precipitation to give the fractional accumulation
380 of precipitation as a function of pentad. Thus, the annual cycle behavior is evaluated irrespective of whether a model
381 has a dry or wet bias. Except for the Gulf of Guinea, the onset and decay of monsoon occur for a fractional
382 accumulation of 0.2 and 0.8, respectively. Between these fractional accumulations, the accumulation of precipitation
383 is nearly linear as the monsoon season progresses. Comparison of the simulated and observed onset, duration, and
384 decay are presented in terms of the difference in the pentad index obtained from the model and observations (i.e.,
385 model minus observations). Therefore, negative values indicate that the onset or decay in the model occurs earlier
386 than in observations, while positive values indicate the opposite. For duration, negative values indicate that for the
387 model it takes fewer pentads to progress from onset to decay compared to observations (i.e., the simulated monsoon
388 period is too short), while positive values indicate the opposite.

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

396

397 **3.6 Cloud feedback and mean-state**

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

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

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

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

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

436

437 **3.7 Precipitation**

438 Recognizing the importance of accurately simulating precipitation in ESMs and a lack of objective and
439 systematic benchmarking for it, and motivated by discussions with WGNE and WGCM working groups of WCRP,
440 the DOE has initiated an effort to establish a pathway to help modelers gauge improvement (U.S. DOE, 2020). The
441 2019 DOE workshop “Benchmarking Simulated Precipitation in Earth System Models” generated two sets of
442 precipitation metrics: *baseline* and *exploratory* metrics (Pendergrass et al., 2020). In the PMP, we have focused on
443 implementing the *baseline* metrics for benchmarking simulated precipitation. In parallel, a set of *exploratory* metrics
444 that could be added to metrics suites including PMP in the future was illustrated by Leung et al. (2022) to extend the
445 evaluation scope to include process-oriented and phenomena-based diagnostics and metrics.

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

452 of the findings here. The precipitation variability metric measures forced (diurnal and annual cycles) and internal
453 variability across timescales (subdaily, synoptic, subseasonal, seasonal, and interannual) in a framework based on
454 power spectra of 3-hourly total and anomaly precipitation. Overall, CMIP5 and CMIP6 models underestimate the
455 internal variability, which is more pronounced in the higher frequency variability, while they overestimate the forced
456 variability (Fig. 10). For the diurnal cycle, PMP includes metrics from Covey et al. (2016). Additionally, the intensity
457 and distribution of precipitation are assessed following Ahn et al. (2023). Extreme daily precipitation indices and their
458 20-year return values are calculated using a non-stationary Generalized Extreme Value statistical method. From the
459 CMIP5 and CMIP6 historical simulations we evaluate model performance of these indices and their return values in
460 comparison with gridded land-based daily observations. Using this approach, Wehner et al. (2020) found that at
461 models' standard resolutions, no meaningful differences were found between the two generations of CMIP models.
462 Wehner et al. (2021) extended the evaluations of simulated extreme precipitation to seasonal 3-hourly precipitation
463 extremes produced by available HighResMIP models and concluded that the improvement is minimal with the models'
464 increased spatial resolutions. They also noted that the order of operations of regridding and calculating extremes
465 affects the ability of models to reproduce observations. Drought metrics developed by Xue and Ullrich (2021) are not
466 implemented in PMP directly, but are wrapped by the Coordinated Model Evaluation Capabilities (CMEC; Ordonez
467 et al. 2021), which is a parallel framework for supporting community-developed evaluation packages. Together, these
468 metrics provide a streamlined workflow for running the entire baseline metrics via the PMP and CMEC that is ready
469 for use by operational centers and in the CMIP7.

470

471 *3.8 Relating metrics to underlying diagnostics*

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

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

491

492 **4 Model Benchmarking**

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

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

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

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

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

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

545

546 **5 Discussion**

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

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

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

581 In addition to the scientific challenges associated with diversifying objective summaries of model
582 performance, there is potential to leverage rapidly evolving technologies, including new open-source tools and
583 methods available to scientists. We expect that the ongoing PMP code modernization effort to fully adapt the xCDAT
584 and xarray will facilitate greater community involvement. As the PMP evolves with these technologies we will
585 continue to maintain rigor in the calculation of statistics for the PMP metrics, for example by incorporating the latest
586 advancements in the field. A prominent example in the objective comparison of models and observations involves the
587 methodology of horizontal interpolation, and in future versions of the PMP we are planning a more stringent
588 conservation method (Taylor, 2023). To improve the clarity of key messages from multivariate PMP metrics data, we
589 will consider implementing the advances in high-dimensional data visualization, e.g., the circular plot discussed in
590 Lee et al. (2018b) and variations of Parallel Coordinate Plots proposed in this paper and by Hassan et al. (2019) and
591 Lu et al. (2020).

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

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

606

607 **6 Summary and Conclusion**

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

623 We expect that the PMP will continue to play a crucial role in benchmarking ESMs. Improvements in the
624 PMP, along with progress in interconnected MIP community projects, will greatly contribute to advancing the
625 evaluation of ESMs including in connection to the community efforts (e.g., the CMIP Benchmarking Task Team).
626 Enhancements in version control and transparency within obs4MIPs are set to enhance the provenance and
627 reproducibility of PMP results, thereby strengthening the foundation for rigorous and repeatable performance
628 benchmarking. The PMP's collaboration with the CMIP Forcing Task Team, through the Input4MIPs (Durack et al.,
629 2018) and the CMIP6Plus projects, will further expand the utility of performance metrics in identifying problems
630 associated with the forcing dataset and their application and use in reproducing the observed record of historical
631 climate. Furthermore, as ESMs advance towards more operationalized configurations to meet the demands of decision-
632 making processes (Jakob et al., 2023), the PMP holds significant potential to provide interoperable ESM evaluation
633 and benchmarking capabilities to the community.

634 **Appendix A: Table of acronyms**

635

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
EWR	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
NPGO	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

637 **Code and Data Availability**

638 The source code of the PMP (Lee et al., 2023b) is available as an open-source Python package:
639 https://github.com/PCMDI/pcmdi_metrics (last access: 21 February 2024) with all released versions archived on
640 Zenodo DOI: <https://doi.org/10.5281/zenodo.592790> (last access: 21 February 2024). The online documentation is
641 available at http://pcmdi.github.io/pcmdi_metrics (last access: 21 February 2024). The PMP results database (Lee et
642 al., 2023a) that includes calculated metrics is available on the GitHub repository at
643 https://github.com/PCMDI/pcmdi_metrics_results_archive (last access: 21 February 2024) with versions archived on
644 Zenodo DOI: <https://doi.org/10.5281/zenodo.10181201>. PMP's installation process is streamlined using the *Anaconda*
645 distribution and the *conda-forge* channel (https://anaconda.org/conda-forge/pcmdi_metrics, last access: 21 February
646 2024). The installation instructions are available at http://pcmdi.github.io/pcmdi_metrics/install.html (last access: 21
647 February 2024). The interactive visualizations of the PMP results are available on the PCMDI website at
648 <https://pcmdi.llnl.gov/metrics> (last access: 21 November 2023). The CMIP5 and CMIP6 model outputs and obs4MIPs
649 datasets used in this paper are available via the Earth System Grid Federation at <https://esgf-node.llnl.gov/> (last access:
650 21 February 2024).

651

652 **Author Contributions**

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

656

657 **Competing interests**

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

660

661 **Acknowledgment**

662 We acknowledge the World Climate Research Programme, which, through its Working Group on Coupled Modeling,
663 coordinated and promoted CMIP6. We thank the climate modeling groups for producing and making available their
664 model output, the Earth System Grid Federation (ESGF) for archiving the data and providing access, and the multiple
665 funding agencies that support CMIP6 and ESGF. This work is performed under the auspices of the U.S. DOE by
666 Lawrence Livermore National Laboratory (LLNL) under Contract No. DE-AC52-07NA27344. Efforts of JL, PJG,
667 MA, AO, PU, KET, PD, CB, MDZ, LC, and BD were supported by the Regional and Global Model Analysis (RGMA)
668 program of the U.S. Department of Energy (DOE) Office of Science (OS), Biological and Environmental Research
669 (BER) program. MFW was supported by the Director, OS, BER of the U.S. DOE through the RGMA program under
670 Contract No. DE340AC02-05CH11231. AGP was supported by U.S. DOE through BER RGMA through Award
671 Number DE-SC0022070 and via National Science Foundation (NSF) IA 1947282, and by National Center for
672 Atmospheric Research (NCAR), which is a major facility sponsored by the NSF under Cooperative Agreement No.
673 1852977. YYP and EG were supported by the Agence Nationale de la Recherche ARISE project, under Grant ANR-

674 18-CE01-0012, and the Belmont project GOTHAM, under Grant ANR-15-JCLI-0004-01, the European
675 Commission’s H2020 Programme “Infrastructure for the European Network for Earth System Modelling Phase 3 (IS-
676 ENES3)” project under Grant Agreement 824084. DK was supported by the New Faculty Startup Fund from Seoul
677 National University and the KMA R&D program (KMI2022-01313). The authors thank Program Manager Renu
678 Joseph of the U.S. DOE for the support and advocacy for the Program for Climate Model Diagnosis and
679 Intercomparison (PCMDI) project and the PMP. We thank Stephen Klein for his leadership for the PCMDI project
680 from 2019 to 2022. We acknowledge contributions from our LLNL colleagues, Lina Muryanto and Zeshawn Shaheen
681 (Now at Google LLC) during the early stage of the PMP, and Sasha Ames, Jeff Painter, Chris Mauzey, and Stephen
682 Po-Chedley for the PCMDI’s CMIP database management. The authors also thank Liping Zhang for her comments
683 during GFDL’s internal review process.

684

685 **References**

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

690 Ahn, M.-S., Kim, D. H., Sperber, K. R., Kang, I.-S., Maloney, E. D., Waliser, D. E., and Hendon, H. H.: MJO
691 simulation in CMIP5 climate models: MJO skill metrics and process-oriented diagnosis, *Climate Dynamics*,
692 49, 4023–4045, <https://doi.org/10.1007/s00382-017-3558-4>, 2017.

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

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

699 Arcodia, M., Barnes, E. A., Mayer, K., Lee, J., Ordonez, A., and Ahn, M.-S.: Assessing decadal variability of
700 subseasonal forecasts of opportunity using explainable AI, *Environmental Research*,
701 <https://doi.org/10.1088/2752-5295/aced60>, 2023.

702 Ashfaq, M., Rastogi, D., Kitson, J., Abid, M. A., and Kao, S.-C.: Evaluation of CMIP6 GCMs over the CONUS for
703 downscaling studies, *Journal of Geophysical Research: Atmospheres*, 127, e2022JD036659.
704 <https://doi.org/10.1029/2022JD036659>, 2022.

705 Bayr, T., Wengel, C., Latif, M., Dommenges, D., Lübbecke, J., and Park, W.: Error compensation of ENSO
706 atmospheric feedbacks in climate models and its influence on simulated ENSO dynamics, *Climate Dynamics*,
707 53, 155–172, <https://doi.org/10.1007/s00382-018-4575-7>, 2019.

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

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

712 Boucher, O., Servonnat, J., Albright, A. L., Aumont, O., Balkanski, Y., Bastrikov, V., Bekki, S., Bonnet, R., Bony,
713 S., Bopp, L., Braconnot, P., Brockmann, P., Cadule, P., Caubel, A., Cheruy, F., Codron, F., Cozic, A., Cugnet,
714 D., D’Andrea, F., Davini, P., De Lavergne, C., Denvil, S., Deshayes, J., Devilliers, M., Ducharne, A.,
715 Dufresne, J. L., Dupont, E., Ethé, C., Fairhead, L., Falletti, L., Flavoni, S., Foujols, M. A., Gardoll, S.,
716 Gastineau, G., Ghattas, J., Grandpeix, J. Y., Guenet, B., Lionel, E. G., Guilyardi, E., Guimberteau, M.,
717 Hauglustaine, D., Hourdin, F., Idelkadi, A., Joussaume, S., Kageyama, M., Khodri, M., Krinner, G., Lebas,
718 N., Levvasseur, G., Lévy, C., Li, L., Lott, F., Lurton, T., Luysaert, S., Madec, G., Madeleine, J.-B.,
719 Maignan, F., Marchand, M., Marti, O., Mellul, L., Meurdesoif, Y., Mignot, J., Musat, I., Ottlé, C., Peylin, P.,
720 Planton, Y., Polcher, J., Rio, C., Rochetin, N., Rousset, C., Sepulchre, P., Sima, A., Swingedouw, D.,
721 Thiéblemont, R., Traore, A. K., Vancoppenolle, M., Vial, J., Vialard, J., Viovy, N., and Vuichard, N.:
722 Presentation and evaluation of the IPSL-CM6A-LR Climate Model, *Journal of Advances in Modeling Earth
723 Systems*, 12, <https://doi.org/10.1029/2019ms002010>, 2020.

724 Caldwell, P., Mametjanov, A., Tang, Q., Van Roekel, L., Golaz, J.-C., Lin, W., Bader, D. C., Keen, N. D., Feng, Y.,
725 Jacob, R., Maltrud, M., Roberts, A., Taylor, M. A., Veneziani, M., Wang, H., Wolfe, J. D., Balaguru, K.,
726 Cameron-Smith, P. J., Dong, L., Klein, S. A., Leung, L. R., Li, H., Li, Q., Liu, X., Neale, R., Pinheiro, M.
727 C., Qian, Y., Ullrich, P. A., Xie, S., Yang, Y., Zhang, Y., Zhang, K., and Zhou, T.: The DOE E3SM Coupled
728 Model Version 1: description and results at high resolution, *Journal of Advances in Modeling Earth Systems*,
729 11, 4095–4146, <https://doi.org/10.1029/2019ms001870>, 2019.

730 Chen, H.-C., Jin, F.-F., Zhao, S., Wittenberg, A. T., and Xie, S.: ENSO dynamics in the E3SM-1-0, CESM2, and
731 GFDL-CM4 climate models, *Journal of Climate*, 34, 9365–9384, <https://doi.org/10.1175/JCLI-D-21-0355.1>,
732 2021.

733 Collier, N., Hoffman, F. M., Lawrence, D. M., Keppel-Aleks, G., Koven, C. D., Riley, W. J., Mu, M., and Randerson,
734 J. T.: The International Land Model Benchmarking (ILAMB) System: Design, theory, and implementation,
735 *Journal of Advances in Modeling Earth Systems*, 10, 2731–2754, <https://doi.org/10.1029/2018ms001354>,
736 2018.

737 Covey, C., AchutaRao, K., Cubasch, U., Jones, P., Lambert, S. J., Mann, M., Phillips, T. J., and Taylor, K. E.: An
738 overview of results from the Coupled Model Intercomparison Project, *Global and Planetary Change*, 37, 103–
739 133, [https://doi.org/10.1016/s0921-8181\(02\)00193-5](https://doi.org/10.1016/s0921-8181(02)00193-5), 2003.

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

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

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

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

748 Dawson, A.: eofs: A Library for EOF Analysis of Meteorological, Oceanographic, and Climate Data, *Journal of Open*
749 *Research Software (JORS)*, 4, e14, <https://doi.org/10.5334/jors.122>, 2016.

750 Dee, D. P., Uppala, S. M., Simmons, A. J., Berrisford, P., Poli, P., Kobayashi, S., Andrae, U., Balmaseda, M. A.,
751 Balsamo, G., Bauer, P., Bechtold, P., Beljaars, A. C. M., van de Berg, L., Bidlot, J., Bormann, N., Delsol,
752 C., Dragani, R., Fuentes, M., Geer, A. J., Haimberger, L., Healy, S. B., Hersbach, H., Hólm, E. V., Isaksen,
753 L., Kållberg, P., Köhler, M., Matricardi, M., McNally, A. P., Monge-Sanz, B. M., Morcrette, J.-J., Park, B.-
754 K., Peubey, C., de Rosnay, P., Tavolato, C., Thépaut, J.-N., and Vitart, F.: The ERA-Interim reanalysis:
755 Configuration and performance of the data assimilation system, *Quarterly Journal of the Royal*
756 *Meteorological Society*, 137, 553–597, <https://doi.org/10.1002/qj.828>, 2011

757 Deser, C. and Phillips, A. S.: Defining the internal component of Atlantic multidecadal variability in a changing
758 climate, *Geophysical Research Letters*, 48, <https://doi.org/10.1029/2021gl095023>, 2021.

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

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

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

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

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

780 Eyring, V., Cox, P. M., Flato, G. M., Gleckler, P. J., Abramowitz, G., Caldwell, P., Collins, W. D., Gier, B. K., Hall,
781 A., Hoffman, F. M., Hurtt, G. C., Jahn, A., Jones, C. D., Klein, S. A., Krasting, J. P., Kwiatkowski, L.,
782 Lorenz, R., Maloney, E. D., Meehl, G. A., Pendergrass, A. G., Pincus, R., Ruane, A. C., Russell, J. L.,

783 Sanderson, B. M., Santer, B. D., Sherwood, S. C., Simpson, I. R., Stouffer, R. J., and Williamson, M. S.:
784 Taking climate model evaluation to the next level, *Nature Climate Change*, 9, 102–110,
785 <https://doi.org/10.1038/s41558-018-0355-y>, 2019.

786 Eyring, V., Bock, L., Lauer, A., Righi, M., Schlund, M., Andela, B., Arnone, E., Bellprat, O., Brötz, B., Caron, L.-P.,
787 Carvalhais, N., Cionni, I., Cortesi, N., Crezee, B., Davin, E. L., Davini, P., Debeire, K., De Mora, L., Deser,
788 C., Docquier, D., Earnshaw, P., Ehbrecht, C., Gier, B. K., Gonzalez-Reviriego, N., Goodman, P. J.,
789 Hagemann, S., Hardiman, S. C., Hassler, B., Hunter, A., Kadow, C., Kindermann, S., Koirala, S., Koldunov,
790 N., Lejeune, Q., Lembo, V., Lovato, T., Lucarini, V., Massonnet, F., Müller, B., Pandde, A., Pérez-Zanón,
791 N., Phillips, A. S., Predoi, V., Russell, J. L., Sellar, A., Serva, F., Stacke, T., Swaminathan, R., Torralba, V.,
792 Vegas-Regidor, J., Von Hardenberg, J., Weigel, K., and Zimmermann, K.: Earth System Model Evaluation
793 Tool (ESMValTool) v2.0 – an extended set of large-scale diagnostics for quasi-operational and
794 comprehensive evaluation of Earth system models in CMIP, *Geoscientific Model Development*, 13, 3383–
795 3438, <https://doi.org/10.5194/gmd-13-3383-2020>, 2020.

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

801 Fasullo, J. T.: Evaluating simulated climate patterns from the CMIP archives using satellite and reanalysis datasets
802 using the Climate Model Assessment Tool (CMATv1), *Geoscientific Model Development*, 13, 3627–3642,
803 <https://doi.org/10.5194/gmd-13-3627-2020>, 2020.

804 Fasullo, J. T., Phillips, A. S., and Deser, C.: Evaluation of leading modes of climate variability in the CMIP archives,
805 *Journal of Climate*, 33, 5527–5545, <https://doi.org/10.1175/jcli-d-19-1024.1>, 2020.

806 Ferraro, R., Waliser, D. E., Gleckler, P. J., Taylor, K. E., and Eyring, V.: Evolving OBS4MIPS to support Phase 6 of
807 the Coupled Model Intercomparison Project (CMIP6), *Bulletin of the American Meteorological Society*,
808 <https://doi.org/10.1175/bams-d-14-00216.1>, 2015.

809 Flato, G., Marotzke, J., Abiodun, B., Braconnot, P., Chou, S.C., Collins, W., Cox, P., Driouech, F., Emori, S., Eyring,
810 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.

813 Fu, W., Moore, J. K., Primeau, F., Collier, N., Ogunro, O. O., Hoffman, F. M. and Randerson, J. T.: Evaluation of
814 ocean biogeochemistry and carbon cycling in CMIP earth system models with the international ocean model
815 benchmarking (IOMB) software System. *Journal of Geophysical Research: Oceans*, 127, e2022JC018965,
816 <https://doi.org/10.1029/2022JC018965>, 2022.

817 Gates, W.L.: AN AMS continuing series: Global CHANGE–AMIP: The Atmospheric Model Intercomparison Project,
818 *Bulletin of the American Meteorological Society*, 73, 1962-1970, 1992.

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

821 Gates, W.L., Boyle, J.S., Covey, C., Dease, C.G., Doutriaux, C.M., Drach, R.S., Fiorino, M., Gleckler, P.J., Hnilo,
822 J.J., Marlais, S.M. and Phillips, T.J.: An overview of the results of the Atmospheric Model Intercomparison
823 Project (AMIP I). *Bulletin of the American Meteorological Society*, 80, 29-56, 1999.

824 Gleckler, P. J., Taylor, K. E., and Doutriaux, C.: Performance metrics for climate models, *Journal of Geophysical*
825 *Research*, 113, <https://doi.org/10.1029/2007jd008972>, 2008.

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

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

831 Golaz, J.-C., Caldwell, P., Van Roekel, L., Petersen, M. R., Tang, Q., Wolfe, J. D., Abeshu, G. W., Anantharaj, V.,
832 Asay-Davis, X., Bader, D. C., Baldwin, S., Bisht, G., Bogenschutz, P., Branstetter, M. L., Brunke, M. A.,
833 Brus, S., Burrows, S. M., Cameron-Smith, P. J., Donahue, A. S., Deakin, M., Easter, R. C., Evans, K. J.,
834 Feng, Y., Flanner, M., Foucar, J. G., Fyke, J., Griffin, B. M., Hannay, C., Harrop, B. E., Hoffman, M. J.,
835 Hunke, E., Jacob, R., Jacobsen, D. W., Jeffery, N., Jones, P. W., Keen, N. D., Klein, S. A., Larson, V. E.,
836 Leung, L. R., Li, H. Y., Lin, W., Lipscomb, W. H., Lun, P., Mahajan, S., Maltrud, M., Mamatjanov, A.,
837 McClean, J. L., McCoy, R., Neale, R., Price, S., Qian, Y., Rasch, P. J., Eyre, J. E. J. R., Riley, W. J., Ringler,
838 T. D., Roberts, A., Roesler, E. L., Salinger, A. G., Shaheen, Z., Shi, X., Singh, B., Tang, J., Taylor, M. A.,
839 Thornton, P. E., Turner, A. K., Veneziani, M., Wan, H., Wang, H., Wang, S., Williams, D. N., Wolfram, P.
840 J., Worley, P. H., Xie, S., Yang, Y., Yoon, J., Zelinka, M. D., Zender, C. S., Zeng, X., Zhang, C., Zhang, K.,
841 Zhang, Y., Zheng, X., Zhou, T., and Zhu, Q.: The DOE E3SM Coupled Model Version 1: Overview and
842 evaluation at standard resolution, *Journal of Advances in Modeling Earth Systems*, 11, 2089–2129,
843 <https://doi.org/10.1029/2018ms001603>, 2019.

844 Goldenson, N., Leung, L. R., Mearns, L. O., Pierce, D. W., Reed, K. A., Simpson, I. R., Ullrich, P., Krantz, W., Hall,
845 A., Jones, A. and Rahimi, S.: Use-Inspired, Process-Oriented GCM Selection: Prioritizing Models for
846 Regional Dynamical Downscaling, *Bulletin of the American Meteorological Society*, E1619–E1629,
847 <https://doi.org/10.1175/BAMS-D-23-0100.1>, 2023.

848 Guilyardi, E., Wittenberg, A., Fedorov, A., Collins, M., Wang, C., Capotondi, A., Van Oldenborgh, G.J. and
849 Stockdale, T.: Understanding El Niño in ocean–atmosphere general circulation models: Progress and
850 challenges, *Bulletin of the American Meteorological Society*, 90, 325-340,
851 <https://doi.org/10.1175/2008BAMS2387.1>, 2009.

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

855 Gutowski Jr., W. J., Giorgi, F., Timbal, B., Frigon, A., Jacob, D., Kang, H.-S., Raghavan, K., Lee, B., Lennard, C.,
856 Nikulin, G., O'Rourke, E., Rixen, M., Solman, S., Stephenson, T., and Tangang, F.: WCRP Coordinated
857 Regional Downscaling EXperiment (CORDEX): a diagnostic MIP for CMIP6, *Geosci. Model Dev.*, 9, 4087–
858 4095, <https://doi.org/10.5194/gmd-9-4087-2016>, 2016.

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

864 Hannah, W. M., Bradley, A. M., Guba, O., Tang, Q., Golaz, J.-C., and Wolfe, J. D.: Separating physics and dynamics
865 grids for improved computational efficiency in spectral element Earth system models, *Journal of Advances*
866 *in Modeling Earth Systems*, 13, <https://doi.org/10.1029/2020ms002419>, 2021.

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

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

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

875 Held, I. M., Guo, H., Adcroft, A., Dunne, J. P., Horowitz, L. W., Krasting, J., Shevliakova, E., Winton, M., Zhao, M.,
876 Bushuk, M., Wittenberg, A. T., and coauthors: Structure and performance of GFDL's CM4. 0 climate model,
877 *Journal of Advances in Modeling Earth Systems*, 11, 3691-3727, <https://doi.org/10.1029/2019MS001829>,
878 2019.

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

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

884 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., Nicolas, J., Peubey, C., Radu, R.,
885 Schepers, D. and coauthors: The ERA5 global reanalysis. *Quarterly Journal of the Royal Meteorological*
886 *Society*, 146, 1999-2049, <https://doi.org/10.1002/qj.3803>, 2020.

887 Hintze, J. L., and Nelson, R. D.: Violin plots: A box plot-density trace synergism, *The American Statistician*,
888 52, 181–184, <https://doi.org/10.1080/00031305.1998.10480559>, 1998.

889 Hoyer, S. and Hamman, J.: xarray: N-D labeled Arrays and Datasets in Python, *Journal of Open Research Software*,
890 5, 10. <https://doi.org/10.5334/jors.148>, 2017.

891 Huffman, G. J., Adler, R. F., Morrissey, M. M., Bolvin, D. T., Curtis, S., Joyce, R., McGavock, B. and Susskind, J.:
892 Global precipitation at one-degree daily resolution from multisatellite observations, *Journal of*
893 *hydrometeorology*, 2, 36-50, 2001.

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

897 Inselberg, A.: Multidimensional detective, in: *Proceedings of IEEE Symposium on Information Visualization*, 100–
898 107, <https://doi.org/10.1109/INFVIS.1997.636793>, 1997.

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

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

904 Johansson, J. and Forsell, C.: Evaluation of parallel coordinates: Overview, categorization and guidelines for future
905 research, *IEEE Transactions on Visualization and Computer Graphics*, 22, 579-588,
906 <https://doi.org/10.1109/TVCG.2015.2466992>, 2016.

907 Jakob, C., Gettelman, A. and Pitman, A.: The need to operationalize climate modelling, *Nat. Clim. Chang.* 13, 1158–
908 1160, <https://doi.org/10.1038/s41558-023-01849-4>, 2023.

909 Joyce, R. J., Janowiak, J. E., Arkin, P. A. and Xie, P.: CMORPH: A method that produces global precipitation
910 estimates from passive microwave and infrared data at high spatial and temporal resolution, *Journal of*
911 *hydrometeorology*, 5, 487-503, 2004.

912 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
913 in CESM2 ensemble simulation, *Geophysical Research Letters*, 47, <https://doi.org/10.1029/2020gl089824>,
914 2020.

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

919 Kim, H., Caron, J. M., Richter, J. H. and Simpson, I. R.: The lack of QBO-MJO connection in CMIP6 models,
920 *Geophysical Research Letters*, 47, e2020GL087295, <https://doi.org/10.1029/2020GL087295>, 2020.

921 Klein, S. A., Zhang, Y., Zelinka, M. D., Pincus, R., Boyle, J., and Gleckler, P. J.: Are climate model simulations of
922 clouds improving? An evaluation using the ISCCP simulator, *Journal of Geophysical Research:*
923 *Atmospheres*, 118, 1329–1342, <https://doi.org/10.1002/jgrd.50141>, 2013.

924 Klingaman, N. P., Martin, G., and Moise, A.: ASoP (v1.0): a set of methods for analyzing scales of precipitation in
925 general circulation models, *Geoscientific Model Development*, 10, 57–83, [https://doi.org/10.5194/gmd-10-](https://doi.org/10.5194/gmd-10-57-2017)
926 [57-2017](https://doi.org/10.5194/gmd-10-57-2017), 2017.

927 Knutti, R.: The end of model democracy? *Climatic Change*, 102, 395–404, <https://doi.org/10.1007/s10584-010-9800->
928 2, 2010.

929 Knutti, R., Sedláček, J., Sanderson, B. M., Lorenz, R., Fischer, E. M., and Eyring, V.: A climate model projection
930 weighting scheme accounting for performance and interdependence, *Geophysical Research Letters*,
931 <https://doi.org/10.1002/2016gl072012>, 2017.

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

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

940 Lambert, S. J. and Boer, G. J.: CMIP1 evaluation and intercomparison of coupled climate models, *Climate Dynamics*,
941 17, 83-106, <https://doi.org/10.1007/PL00013736>, 2001.

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

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

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

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

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

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

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

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

967 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.,
 968 Klingaman, N. P. Kuo, Y.-H., Lee, R. W., Martinez-Villalobos, C., Vishnu S., Priestley, M. D. K., Tao, C.,
 969 and Zhou, Y.: Exploratory precipitation metrics: Spatiotemporal characteristics, process-oriented, and
 970 phenomena-based, *Journal of Climate*, 35, <https://doi.org/10.1175/JCLI-D-21-0590.1>, 3659-3686, 2022.

971 Lin, J.-P., Kiladis, G. N., Mapes, B. E., Weickmann, K. M., Sperber, K. R., Lin, W., Wheeler, M. C., Schubert, S. D.,
 972 Del Genio, A. D., Donner, L. J., Emori, S., Gu er emy, J.-F., Hourdin, F., Rasch, P. J., Roeckner, E., and
 973 Scinocca, J.: Tropical intraseasonal variability in 14 IPCC AR4 climate Models. Part I: Convective Signals,
 974 *Journal of Climate*, 19, 2665–2690, <https://doi.org/10.1175/jcli3735.1>, 2006.

975 Lin, Y., Huang, X., Liang, Y., Qin, Y., Xu, S., Huang, W., Xu, F., Liu, L., Wang, Y., Peng, Y. and Wang, L.:
 976 Community integrated earth system model (CIESM): Description and evaluation, *Journal of Advances in*
 977 *Modeling Earth Systems*, 12, <https://doi.org/10.1029/2019ms002036>, 2020.

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

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

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

987 Madden, R. A. and Julian, P.: Detection of a 40–50 day oscillation in the zonal wind in the Tropical Pacific, *Journal*
 988 *of the Atmospheric Sciences*, 28, 702–708, [https://doi.org/10.1175/1520-0469\(1971\)028](https://doi.org/10.1175/1520-0469(1971)028), 1971.

989 Madden, R. A. and Julian, P.: Description of Global-Scale Circulation Cells in the Tropics with a 40–50 Day Period,
 990 *Journal of the Atmospheric Sciences*, 29, 1109–1123, [https://doi.org/10.1175/1520-0469\(1972\)029](https://doi.org/10.1175/1520-0469(1972)029), 1972.

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

993 Martin, G. M., Klingaman, N. P., and Moise, A. F.: Connecting spatial and temporal scales of tropical precipitation in
 994 observations and the MetUM-GA6, *Geoscientific Model Development*, 10, 105–126,
 995 <https://doi.org/10.5194/gmd-10-105-2017>, 2017.

996 Maloney, E. D., Gettelman, A., Ming, Y., Neelin, J. D., Barrie, D., Mariotti, A., Chen, C., Coleman, D., Kuo, Y. H.,
 997 Singh, B., Annamalai, H., Berg, A., Booth, J. F., Camargo, S. J., Dai, A., Gonzalez, A., Hafner, J., Jiang, X.,
 998 Jing, X., Kim, D. H., Kumar, A., Moon, Y., Naud, C. M., Sobel, A. H., Suzuki, K., Wang, F., Wang, J., Wing,
 999 A. A., Xu, X., and Zhao, M.: Process-Oriented evaluation of climate and weather forecasting models, *Bulletin*
 1000 *of the American Meteorological Society*, 100, 1665–1686, <https://doi.org/10.1175/bams-d-18-0042.1>, 2019.

1001 McAvaney, B.J., Covey, C., Joussaume, S., Kattsov, V., Kitoh, A., Ogana, W., Pitman, A.J., Weaver, A.J., Wood,
1002 R.A. and Zhao, Z.C.: Model evaluation. In *Climate Change 2001: The scientific basis. Contribution of WG1*
1003 *to the Third Assessment Report of the IPCC (TAR) 471-523*, Cambridge University Press, 2001.

1004 McPhaden, M. J., Zebiak, S. E., and Glantz, M. H.: ENSO as an integrating concept in Earth Science, *Science*, 314,
1005 1740–1745, <https://doi.org/10.1126/science.1132588>, 2006.

1006 McPhaden, M. J., Santoso, A., Cai, W. (Eds.): *El Niño Southern oscillation in a changing climate*, American
1007 Geophysical Union, USA, 528 pp., ISBN:9781119548126, <https://doi.org/10.1002/9781119548164>, 2020.

1008 Mears, C. A., Smith, D. K., Ricciardulli, L., Wang, J., Huelsing, H., & Wentz, F. J.: Construction and uncertainty
1009 estimation of a satellite-derived total precipitable water data record over the world's oceans, *Earth and Space*
1010 *Science*, 5, 197–210, <https://doi.org/10.1002/2018EA000363>, 2018.

1011 Meehl, G. A., Boer, G. J., Covey, C., Latif, M., and Stouffer, R. J.: The Coupled Model Intercomparison Project
1012 (CMIP), *Bulletin of the American Meteorological Society*, 81, 313–318, 2000.

1013 Meehl, G. A., Boer, G. J., Covey, C., Latif, M., and Stouffer, R. J.: Intercomparison makes for a better climate model,
1014 *Eos, Transactions American Geophysical Union*, 78, 445, <https://doi.org/10.1029/97eo00276>, 1997.

1015 Meehl, G. A., Covey, C., Delworth, T. L., Latif, M., McAvaney, B. J., Mitchell, J. F. B., Stouffer, R. J., and Taylor,
1016 K. E.: THE WCRP CMIP3 Multimodel Dataset: A new era in climate change research, *Bulletin of the*
1017 *American Meteorological Society*, 88, 1383–1394, <https://doi.org/10.1175/bams-88-9-1383>, 2007.

1018 Merrifield, A., Brunner, L., Lorenz, R., Humphrey, V., and Knutti, R.: Climate model Selection by Independence,
1019 Performance, and Spread (ClimSIPS v1.0.1) for regional applications, *Geoscientific Model Development*,
1020 16, 4715–4747, <https://doi.org/10.5194/gmd-16-4715-2023>, 2023.

1021 Neelin, J. D., Krasting, J. P., Radhakrishnan, A., Liptak, J., Jackson, T. J., Ming, Y., Dong, W., Gettelman, A.,
1022 Coleman, D., Maloney, E. D., Wing, A. A., Kuo, Y. H., Ahmed, F., Ullrich, P. A., Bitz, C. M., Neale, R.,
1023 Ordonez, A., and Maroon, E.: Process-oriented diagnostics: principles, practice, community development
1024 and common standards, *Bulletin of the American Meteorological Society*, [https://doi.org/10.1175/bams-d-](https://doi.org/10.1175/bams-d-21-0268.1)
1025 [21-0268.1](https://doi.org/10.1175/bams-d-21-0268.1), 2023.

1026 Nowack, P., Runge, J., Eyring, V., and Haigh, J. D.: Causal networks for climate model evaluation and constrained
1027 projections, *Nature Communications*, 11, <https://doi.org/10.1038/s41467-020-15195-y>, 2020.

1028 Orbe, C., Van Roekel, L., Adames, Á. F., Dezfuli, A., Fasullo, J. T., Gleckler, P. J., Lee, J., Li, W., Nazarenko, L.,
1029 Schmidt, G. A., Sperber, K. R., and Zhao, M.: Representation of modes of variability in six U.S. climate
1030 models, *Journal of Climate*, 33, 7591–7617, <https://doi.org/10.1175/jcli-d-19-0956.1>, 2020.

1031 Ordonez, A. C., Klingaman, N. P., and Martin, G.: Analysing scales of precipitation, OSTI OAI (U.S. Department of
1032 Energy Office of Scientific and Technical Information), <https://doi.org/10.11578/dc.20211029.5>, 2021.

1033 Papalexiou, S. M., Rajulapati, C. R., Clark, M. P., and Lehner, F.: Robustness of CMIP6 historical global mean
1034 temperature simulations: Trends, long-term persistence, autocorrelation, and distributional shape, *Earth's*
1035 *Future*, 8, e2020EF001667, <https://doi.org/10.1029/2020EF001667>, 2020.

1036 Pascoe, C., Lawrence, B. N., Guilyardi, E., Juckes, M., and Taylor, K. E.: Documenting numerical experiments in
1037 support of the Coupled Model Intercomparison Project Phase 6 (CMIP6), *Geosci. Model Dev.*, 13, 2149–
1038 2167, <https://doi.org/10.5194/gmd-13-2149-2020>, 2020.

1039 Pendergrass, A. G., Gleckler, P. J., Leung, L. R., and Jakob, C.: Benchmarking simulated precipitation in earth system
1040 models, *Bulletin of the American Meteorological Society*, 101, E814–E816, <https://doi.org/10.1175/bams-d-19-0318.1>, 2020.

1042 Phillips, A. S., Deser, C., and Fasullo, J. T.: Evaluating modes of variability in climate models, *Eos, Transactions*
1043 *American Geophysical Union*, 95, 453–455, <https://doi.org/10.1002/2014eo490002>, 2014.

1044 Planton, Y., Guilyardi, E., Wittenberg, A. T., Lee, J., Gleckler, P. J., Bayr, T., McGregor, S., McPhaden, M. J., Power,
1045 S. B., Roehrig, R., Vialard, J., and Voldoire, A.: Evaluating Climate Models with the CLIVAR 2020 ENSO
1046 Metrics Package, *Bulletin of the American Meteorological Society*, 102, E193–E217,
1047 <https://doi.org/10.1175/bams-d-19-0337.1>, 2021.

1048 Planton, Y. Y., Lee, J., Wittenberg, A. T., Gleckler, P. J., Guilyardi, E., McGregor, S., and McPhaden, M. J.:
1049 Estimating uncertainty in simulated ENSO statistics, *Journal of Advances in Modeling Earth Systems* (under
1050 review), ESS Open Archive, <https://doi.org/10.22541/essoar.170196744.48068128/v1>, 2023.

1051 Potter, G. L., Bader, D. C., Riches, M., Bamzai, A. and Joseph, R.: Celebrating two decades of the Program for Climate
1052 Model Diagnosis and Intercomparison. *Bulletin of the American Meteorological Society*, 92, 629–631,
1053 <https://doi.org/10.1175/2011BAMS3018.1>, 2011.

1054 Qin, Y., Zelinka, M. D., and Klein, S. A.: On the Correspondence Between Atmosphere-Only and Coupled
1055 Simulations for Radiative Feedbacks and Forcing From CO₂, *Journal of Geophysical Research:*
1056 *Atmospheres*, 127, <https://doi.org/10.1029/2021jd035460>, 2022.

1057 Randall, D.A., Wood, R.A., Bony, S., Colman, R., Fichefet, T., Fyfe, J., Kattsov, V., Pitman, A., Shukla, J., Srinivasan,
1058 J. and Stouffer, R.J.: Climate models and their evaluation. In *Climate change 2007: The physical science*
1059 *basis. Contribution of Working Group I to the Fourth Assessment Report of the IPCC (FAR)*, 589–662,
1060 Cambridge University Press, 2007.

1061 Randerson, J. T., Hoffman, F. M., Thornton, P. E., Mahowald, N. M., Lindsay, K., Lee, Y.-H., Nevison, C. D., Doney,
1062 S. C., Bonan, G. B., Stöckli, R., Covey, C., Running, S. W., and Fung, I.: Systematic assessment of terrestrial
1063 biogeochemistry in coupled climate-carbon models, *Global Change Biology*, 15, 2462–2484,
1064 <https://doi.org/10.1111/j.1365-2486.2009.01912.x>, 2009.

1065 Rasch, P. J., Xie, S., Ma, P.-L., Lin, W., Wang, H., Tang, Q., Burrows, S. M., Caldwell, P., Zhang, K., Easter, R. C.,
1066 Cameron-Smith, P. J., Singh, B., Wan, H., Golaz, J.-C., Harrop, B. E., Roesler, E. L., Bacmeister, J. T.,
1067 Larson, V. E., Evans, K. J., Qian, Y., Taylor, M. A., Leung, L. R., Zhang, Y., Brent, L., Branstetter, M. L.,
1068 Hannay, C., Mahajan, S., Mامتjanov, A., Neale, R., Richter, J. H., Yoon, J.-H., Zender, C. S., Bader, D. C.,
1069 Flanner, M., Foucar, J. G., Jacob, R., Keen, N. D., Klein, S. A., Liu, X., Salinger, A. G., Shrivastava, M., and
1070 Yang, Y.: An overview of the atmospheric component of the Energy Exascale Earth System model, *Journal*
1071 *of Advances in Modeling Earth Systems*, 11, 2377–2411, <https://doi.org/10.1029/2019ms001629>, 2019.

1072 Reichler, T. and Kim, J.: How well do coupled models simulate today's climate?, *Bulletin of the American*
1073 *Meteorological Society*, 89, 303–312, <https://doi.org/10.1175/bams-89-3-303>, 2008.

1074 Righi, M., Andela, B., Eyring, V., Lauer, A., Predoi, V., Schlund, M., Vegas-Regidor, J., Bock, L., Brötz, B., De
1075 Mora, L., Diblen, F., Dreyer, L., Drost, N., Earnshaw, P., Hassler, B., Koldunov, N., Little, B., Tomas, S. L.,
1076 and Zimmermann, K.: Earth System Model Evaluation Tool (ESMValTool) v2.0 – technical overview,
1077 *Geoscientific Model Development*, 13, 1179–1199, <https://doi.org/10.5194/gmd-13-1179-2020>, 2020.

1078 Sanderson, B. M. and Wehner, M. F.: Weighting strategy for the Fourth National Climate Assessment, in: *Climate*
1079 *Science Special Report: Fourth National Climate Assessment, Volume I*, edited by Wuebbles, D. J., Fahey,
1080 D. W., Hibbard, K. A., Dokken, D. J., Stewart, B. C., and Maycock, T.K., U.S. Global Change Research
1081 Program, Washington, DC, USA, 436-442, <https://doi.org/10.7930/J06T0JS3>, 2017.

1082 Sanderson, B. M., Wehner, M., and Knutti, R.: Skill and independence weighting for multi-model assessments,
1083 *Geoscientific Model Development*, 10, 2379–2395, <https://doi.org/10.5194/gmd-10-2379-2017>, 2017.

1084 Sherwood, S. C., Webb, M. J., Annan, J. D., Armour, K. C., Forster, P. M., Hargreaves, J. C., Hegerl, G. C., Klein, S.
1085 A., Marvel, K., Rohling, E. J., Watanabe, M., Andrews, T., Braconnot, P., Bretherton, C. S., Foster, G. L.,
1086 Hausfather, Z., Von Der Heydt, A. S., Knutti, R., Mauritsen, T., Norris, J. R., Proistosescu, C., Rugenstein,
1087 M., Schmidt, G. A., Tokarska, K., and Zelinka, M. D.: An assessment of Earth's climate sensitivity using
1088 multiple lines of evidence, *Reviews of Geophysics*, 58, <https://doi.org/10.1029/2019rg000678>, 2020.

1089 Singh, R., and AchutaRao, K.: Sensitivity of future climate change and uncertainty over India to
1090 performance-based model weighting. *Clim. Change*, <https://doi.org/10.1007/s10584-019-02643-y>, 2020.

1091 Sillmann, J., Kharin, V. V., Zhang, X., Zwiers, F. W., and Bronaugh, D.: Climate extremes indices in the CMIP5
1092 multimodel ensemble: Part 1. Model evaluation in the present climate, *Journal of Geophysical Research:*
1093 *Atmospheres*, 118, 1716–1733, <https://doi.org/10.1002/jgrd.50203>, 2013.

1094 Sperber, K. R.: Madden-Julian variability in NCAR CAM2.0 and CCSM2.0, *Clim Dyn* 23, 259–278,
1095 <https://doi.org/10.1007/s00382-004-0447-4>, 2004.

1096 Sperber, K. R., Annamalai, H., Kang, I.-S., Kitoh, A., Moise, A., Turner, A., Wang, B., and Zhou, T.: The Asian
1097 summer monsoon: an intercomparison of CMIP5 vs. CMIP3 simulation of the late 20th century. *Clim Dyn*
1098 41, 2711-2744, <https://doi.org/10.1007/s00382-012-1607-6>, 2013.

1099 Sperber K. R., Gualdi, S., Legutke, S., Gayler, V.: The Madden–Julian oscillation in ECHAM4 coupled and uncoupled
1100 general circulation models, *Clim Dyn* 25, 117–140, <https://doi.org/10.1007/s00382-005-0026-3>, 2005.

1101 Srivastava, A., Grotjahn, R., and Ullrich, P. A.: Evaluation of historical CMIP6 model simulations of extreme
1102 precipitation over contiguous US regions, *Weather and Climate Extremes*, 29, 100268,
1103 <https://doi.org/10.1016/j.wace.2020.100268>, 2020.

1104 Steed, C. A., Shipman, G., Thornton, P., Ricciuto, D., Erickson, D. and Branstetter, M.: Practical application of parallel
1105 coordinates for climate model analysis, *Procedia Computer Science*, 9, 877-886,
1106 <https://doi.org/10.1016/j.procs.2012.04.094>, 2012.

1107 Stevens, B., Satoh, M., Auger, L., Biercamp, J., Bretherton, C. S., Chen, X., Düben, P., Judt, F., Khairoutdinov, M.,
1108 Klocke, D., Kodama, C., Kornblueh, L., Lin, S.-J., Neumann, P., Putman, W. M., Röber, N., Shibuya, R.,

1109 Vanniere, B., Vidale, P. L., Wedi, N., and Zhou, L.: DYAMOND: the DYNAMICS of the Atmospheric general
1110 circulation Modeled On Non-hydrostatic Domains, *Progress in Earth and Planetary Science*, 6,
1111 <https://doi.org/10.1186/s40645-019-0304-z>, 2019.

1112 Stoner, A. M. K., Hayhoe, K., and Wuebbles, D. J.: Assessing general circulation model simulations of atmospheric
1113 teleconnection patterns, *Journal of Climate*, 22, 4348–4372, <https://doi.org/10.1175/2009jcli2577.1>, 2009.

1114 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,
1115 Y., Lee, J., Lee, J., Kim, J.-S., Marzin, C., and Byun, Y.-H.: Climate change projection in the Twenty-First
1116 Century simulated by NIMS-KMA CMIP6 model based on new GHGs concentration pathways, *Asia-pacific
1117 Journal of Atmospheric Sciences*, 57, 851–862, <https://doi.org/10.1007/s13143-021-00225-6>, 2021.

1118 Tang, Q., Prather, M. J., Hsu, J., Ruiz, D. J., Cameron-Smith, P. J., Xie, S., and Golaz, J.-C.: Evaluation of the
1119 interactive stratospheric ozone (O3v2) module in the E3SM version 1 Earth system model, *Geoscientific
1120 Model Development*, 14, 1219–1236, <https://doi.org/10.5194/gmd-14-1219-2021>, 2021.

1121 Tang, S., Fast, J. D., Zhang, K., Hardin, J. C., Varble, A. C., Shilling, J. E., Mei, F., Zawadowicz, M. A., and Ma, P.-
1122 L.: Earth System Model Aerosol–Cloud Diagnostics (ESMAC Diags) package, version 1: assessing E3SM
1123 aerosol predictions using aircraft, ship, and surface measurements, *Geosci. Model Dev.*, 15, 4055–4076,
1124 <https://doi.org/10.5194/gmd-15-4055-2022>, 2022.

1125 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.:
1126 Earth System Model Aerosol–Cloud Diagnostics (ESMAC Diags) package, version 2: assessing aerosols,
1127 clouds, and aerosol–cloud interactions via field campaign and long-term observations, *Geosci. Model Dev.*,
1128 16, 6355–6376, <https://doi.org/10.5194/gmd-16-6355-2023>, 2023.

1129 Taylor, K. E.: Truly Conserving with Conservative Remapping Methods, *Geosci. Model Dev. Discuss.* [preprint],
1130 <https://doi.org/10.5194/gmd-2023-177>, in review, 2023.

1131 Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, *Journal of Geophysical
1132 Research*, 106, 7183–7192, <https://doi.org/10.1029/2000jd900719>, 2001.

1133 Taylor, K. E., Stouffer, R. J., and Meehl, G. A.: An overview of CMIP5 and the experiment design, *Bulletin of the
1134 American Meteorological Society*, 93, 485–498, <https://doi.org/10.1175/bams-d-11-00094.1>, 2012.

1135 Teixeira, J., Waliser, D. E., Ferraro, R., Gleckler, P. J., Lee, T., and Potter, G. L.: Satellite observations for CMIP5:
1136 The Genesis of OBS4MIPs, *Bulletin of the American Meteorological Society*, 95, 1329–1334,
1137 <https://doi.org/10.1175/bams-d-12-00204.1>, 2014.

1138 Tian, B., and Dong, X.: The Double-ITCZ Bias in CMIP3, CMIP5, and CMIP6 Models Based on Annual Mean
1139 Precipitation, *Geophysical Research Letters*, 47, e2020GL087232, <https://doi.org/10.1029/2020GL087232>,
1140 2020

1141 Ullrich, P. A. and Zarzycki, C. M.: TempestExtremes: a framework for scale-insensitive pointwise feature tracking on
1142 unstructured grids, *Geoscientific Model Development*, 10, 1069–1090, <https://doi.org/10.5194/gmd-10-1069-2017>, 2017.

1144 Ullrich, P. A., Zarzycki, C. M., McClenny, E., Pinheiro, M. C., Stansfield, A. M., and Reed, K. A.: TempestExtremes
1145 v2.1: a community framework for feature detection, tracking, and analysis in large datasets, *Geoscientific*
1146 *Model Development*, 14, 5023–5048, <https://doi.org/10.5194/gmd-14-5023-2021>, 2021.

1147 U.S. Department of Energy (DOE): Benchmarking Simulated Precipitation in Earth System Models Workshop Report,
1148 DOE/SC-0203, U.S. Department of Energy Office of Science, Biological and Environmental Research (BER)
1149 Program. Germantown, Maryland, USA. 2020.

1150 Vo, T., Po-Chedley, P., Boutte, J., Zhang, C., Lee, J., Gleckler, P., Durack, P., Taylor, K., and Golaz, J.-C.: Xarray
1151 Climate Data Analysis Tools (xCDAT): A Python Package for Simple and Robust Analysis of Climate Data,
1152 The 103rd AMS Annual Meeting, Abstract, 2023.

1153 Waliser, D. E., Gleckler, P. J., Ferraro, R., Taylor, K. E., Ames, S., Biard, J., Bosilovich, M. G., Brown, O. B., Chepfer,
1154 H., Cinquini, L., Durack, P. J., Eyring, V., Mathieu, P.-P., Lee, T., Pinnock, S., Potter, G. L., Rixen, M.,
1155 Saunders, R., Schulz, J. B., Thépaut, J.-N., and Tuma, M.: Observations for Model Intercomparison Project
1156 (Obs4MIPs): status for CMIP6, *Geoscientific Model Development*, 13, 2945–2958,
1157 <https://doi.org/10.5194/gmd-13-2945-2020>, 2020.

1158 Waliser, D. E., Sperber, K. R., Hendon, H. H., Kim, D., Maloney, E. D., Wheeler, M. C., Weickmann, K. M., Zhang,
1159 C., Donner, L. J., Gottschalck, J., Higgins, W., Kang, I. S., Legler, D. M., Moncrieff, M. W., Schubert, S. D.,
1160 Stern, W., Vitart, F., Wang, B., Wang, W., and Woolnough, S. J.: MJO Simulation Diagnostics, *Journal of*
1161 *Climate*, 22, 3006–3030, <https://doi.org/10.1175/2008jcli2731.1>, 2009.

1162 Wang, J., Liu, X., Shen, H. W. and Lin, G.: Multi-resolution climate ensemble parameter analysis with nested parallel
1163 coordinates plots, *IEEE Transactions on Visualization and Computer Graphics*, 23, 81-90,
1164 <https://doi.org/10.1109/TVCG.2016.2598830>, 2017.

1165 Wehner, M., Gleckler, P. J., and Lee, J.: Characterization of long period return values of extreme daily temperature
1166 and precipitation in the CMIP6 models: Part 1, model evaluation, *Weather and Climate Extremes*, 30,
1167 100283, <https://doi.org/10.1016/j.wace.2020.100283>, 2020.

1168 Wehner, M., Lee, J., Risser, M. D., Ullrich, P. A., Gleckler, P. J., and Collins, W. D.: Evaluation of extreme sub-daily
1169 precipitation in high-resolution global climate model simulations, *Philosophical Transactions of the Royal*
1170 *Society A*, 379, 20190545, <https://doi.org/10.1098/rsta.2019.0545>, 2021.

1171 Whitehall, K., Mattmann, C., Waliser, D., Kim, J., Goodale, C., Hart, A., Ramirez, P., Zimdars, P., Crichton, D.,
1172 Jenkins, G., Jones, C., Asrar, G., and Hewitson, B.: Building Model Evaluation and Decision Support
1173 Capacity for CORDEX, *WMO Bulletin*, 61, available at:
1174 [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)
1175 [cordex](https://public.wmo.int/en/resources/bulletin/building-model-evaluation-and-decision-support-capacity-cordex) (last access date: 14 September 2023), 2012.

1176 Williams, D. N.: Visualization and analysis tools for ultrascale climate data, *Eos, Transactions American Geophysical*
1177 *Union*, 95, 377–378, <https://doi.org/10.1002/2014eo420002>, 2014.

1178 Williams, D. N., Balaji, V., Cinquini, L., Denvil, S., Duffy, D. Q., Evans, B., Ferraro, R., Hansen, R., Lautenschlager,
1179 M., and Trenham, C.: A global repository for Planet-Sized experiments and observations, *Bulletin of the*
1180 *American Meteorological Society*, 97, 803–816, <https://doi.org/10.1175/bams-d-15-00132.1>, 2016.

1181 Williams, D. N., Doutriaux, C., Drach, R., and McCoy, R.: The Flexible Climate Data Analysis Tools (CDAT) for
1182 Multi-model Climate Simulation Data, IEEE International Conference on Data Mining Workshops, 254–261,
1183 <https://doi.org/10.1109/icdmw.2009.64>, 2009.

1184 Wong, P. C., Shen, H. W., Leung, R., Hagos, S., Lee, T. Y., Tong, X. and Lu, K.: Visual analytics of large-scale
1185 climate model data, in: 2014 IEEE 4th Symposium on Large Data Analysis and Visualization (LDAV), 85-
1186 92, <https://doi.org/10.1109/LDAV.2014.7013208>, 2014.

1187 Xie, P., Joyce, R., Wu, S., Yoo, S.H., Yarosh, Y., Sun, F. and Lin, R.: Reprocessed, bias-corrected CMORPH global
1188 high-resolution precipitation estimates from 1998, *Journal of Hydrometeorology*, 18, 1617-1641, 2017.

1189 Xue, Z. and Ullrich, P. A.: A Comprehensive Intermediate-Term Drought Evaluation System and Evaluation of
1190 Climate Data Products over the Conterminous United States, *Journal of Hydrometeorology*,
1191 <https://doi.org/10.1175/jhm-d-20-0314.1>, 2021.

1192 Young, A. H., Knapp, K. R., Inamdar, A. K., Hankins, W., and Rossow, W. B.: The International Satellite Cloud
1193 Climatology Project H-Series climate data record product, *Earth System Science Data*, 10, 583–593,
1194 <https://doi.org/10.5194/essd-10-583-2018>, 2018.

1195 Zelinka, M. D., Klein, S. A., Qin, Y., and Myers, T. A.: Evaluating climate models’ cloud feedbacks against expert
1196 judgment, *Journal of Geophysical Research: Atmospheres*, 127, <https://doi.org/10.1029/2021jd035198>,
1197 2022.

1198 Zelinka, M. D., Myers, T. A., McCoy, D. T., Po-Chedley, S., Caldwell, P. M., Ceppi, P., Klein, S. A., and Taylor, K.
1199 E.: Causes of higher climate sensitivity in CMIP6 models, *Geophysical Research Letters*, 47,
1200 e2019GL085782, <https://doi.org/10.1029/2019GL085782>, 2020.

1201 Zhang, C., Golaz, J.-C., Forsyth, R., Vo, T., Xie, S., Shaheen, Z., Potter, G. L., Asay-Davis, X. S., Zender, C. S., Lin,
1202 W., Chen, C.-C., Terai, C. R., Mahajan, S., Zhou, T., Balaguru, K., Tang, Q., Tao, C., Zhang, Y.,
1203 Emmenegger, T., Burrows, S., and Ullrich, P. A.: The E3SM Diagnostics Package (E3SM Diags v2.7): a
1204 Python-based diagnostics package for Earth system model evaluation, *Geosci. Model Dev.*, 15, 9031–9056,
1205 <https://doi.org/10.5194/gmd-15-9031-2022>, 2022.

1206 Zhang, C. and Hendon, H. H.: Propagating and standing components of the intraseasonal oscillation in tropical
1207 convection, *Journal of the Atmospheric Sciences*, 54, 741–752, [https://doi.org/10.1175/1520-0469\(1997\)054](https://doi.org/10.1175/1520-0469(1997)054), 1997.

1209 Zhang, C., Xie, S., Klein, S. A., Ma, H. Y., Tang, S., Van Weverberg, K., Morcrette, C. J. and Petch, J.: CAUSES:
1210 Diagnosis of the summertime warm bias in CMIP5 climate models at the ARM Southern Great Plains site.
1211 *Journal of Geophysical Research: Atmospheres*, 123, 2968-2992, <https://doi.org/10.1002/2017JD027200>,
1212 2018.

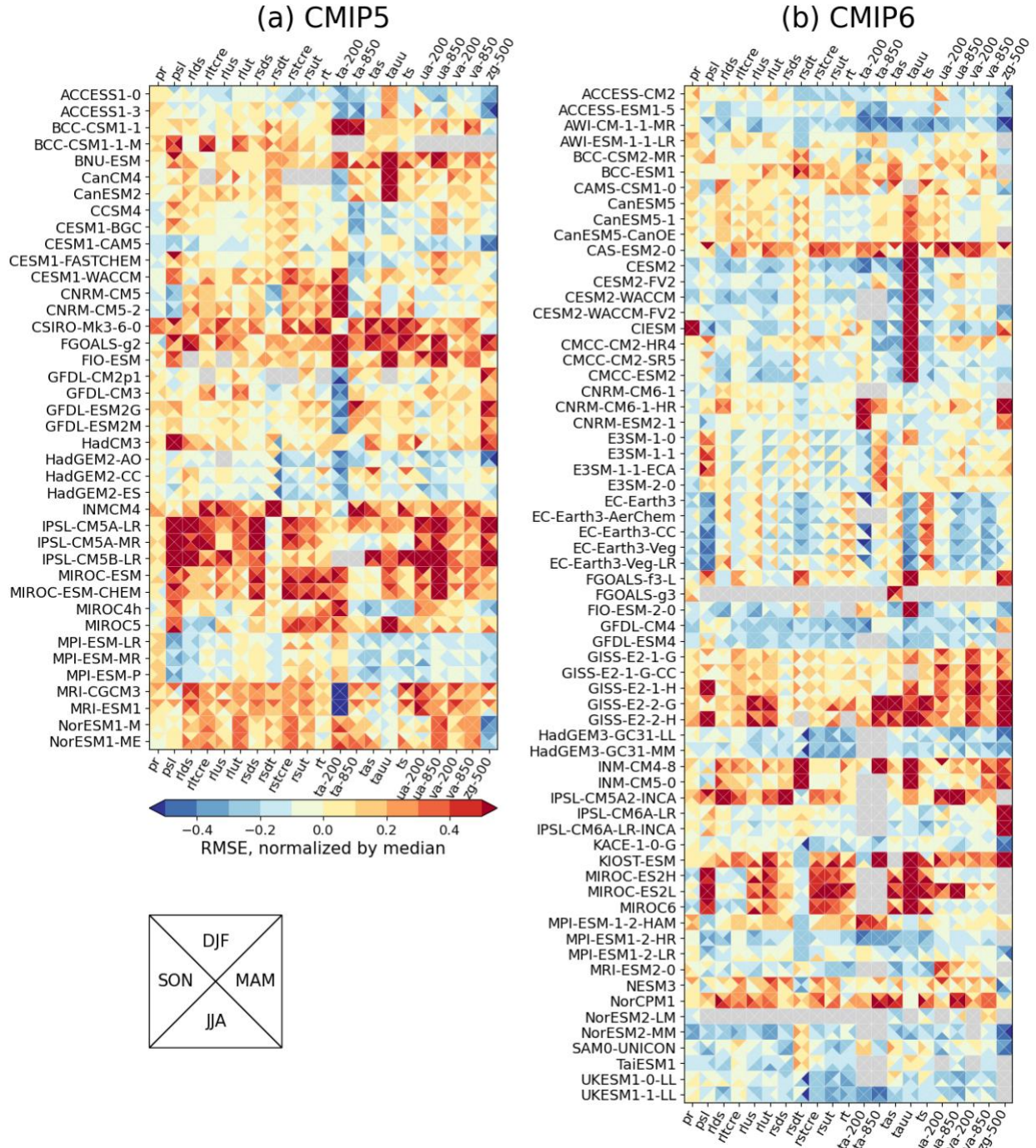
1213 Zhang, C., Xie, S., Tao, C., Tang, S., Emmenegger, T., Neelin, J. D., Schiro, K. A., Lin, W. and Shaheen, Z.: The
1214 ARM data-oriented metrics and diagnostics package for climate models: A new tool for evaluating climate
1215 models with field data, <https://doi.org/10.1175/BAMS-D-19-0282.1>, *Bulletin of the American*
1216 *Meteorological Society*, 101, E1619-E1627, 2020.

1217 Zhao, B., Lin, P., Hu, A., Liu, H., Ding, M., Yu, Z., and Yu, Y.: Uncertainty in Atlantic Multidecadal Oscillation
1218 derived from different observed datasets and their possible causes, *Frontiers in Marine Science*, 9,
1219 <https://doi.org/10.3389/fmars.2022.1007646>, 2022.

1220 Zhao, M., Golaz, J.-C., Held, I. M., Guo, H., Balaji, V., Benson, R., Chen, J. H., Chen, X., Donner, L. J., Dunne, J.,
1221 Dunne, K. A., Durachta, J., Fan, S.-M., Freidenreich, S. M., Garner, S. T., Ginoux, P., Harris, L., Horowitz,
1222 L. W., Krasting, J. P., Langenhorst, A. R., Zhi, L., Lin, P., Lin, S. J., Malyshev, S., Mason, E., Milly, P. C.
1223 D., Ming, Y., Naik, V., Paulot, F., Paynter, D., Phillipps, P. J., Radhakrishnan, A., Ramaswamy, V.,
1224 Robinson, T., Schwarzkopf, D., Seman, C. J., Shevliakova, E., Shen, Z., Shin, H. H., Silvers, L. G., Wilson,
1225 J. R., Winton, M., Wittenberg, A. T., Wyman, B., and Xiang, B.: The GFDL Global Atmosphere and Land
1226 Model AM4.0/LM4.0: 1. Simulation characteristics with prescribed SSTs, *Journal of Advances in Modeling
1227 Earth Systems*, 10, 691–734, <https://doi.org/10.1002/2017ms001208>, 2018.
1228

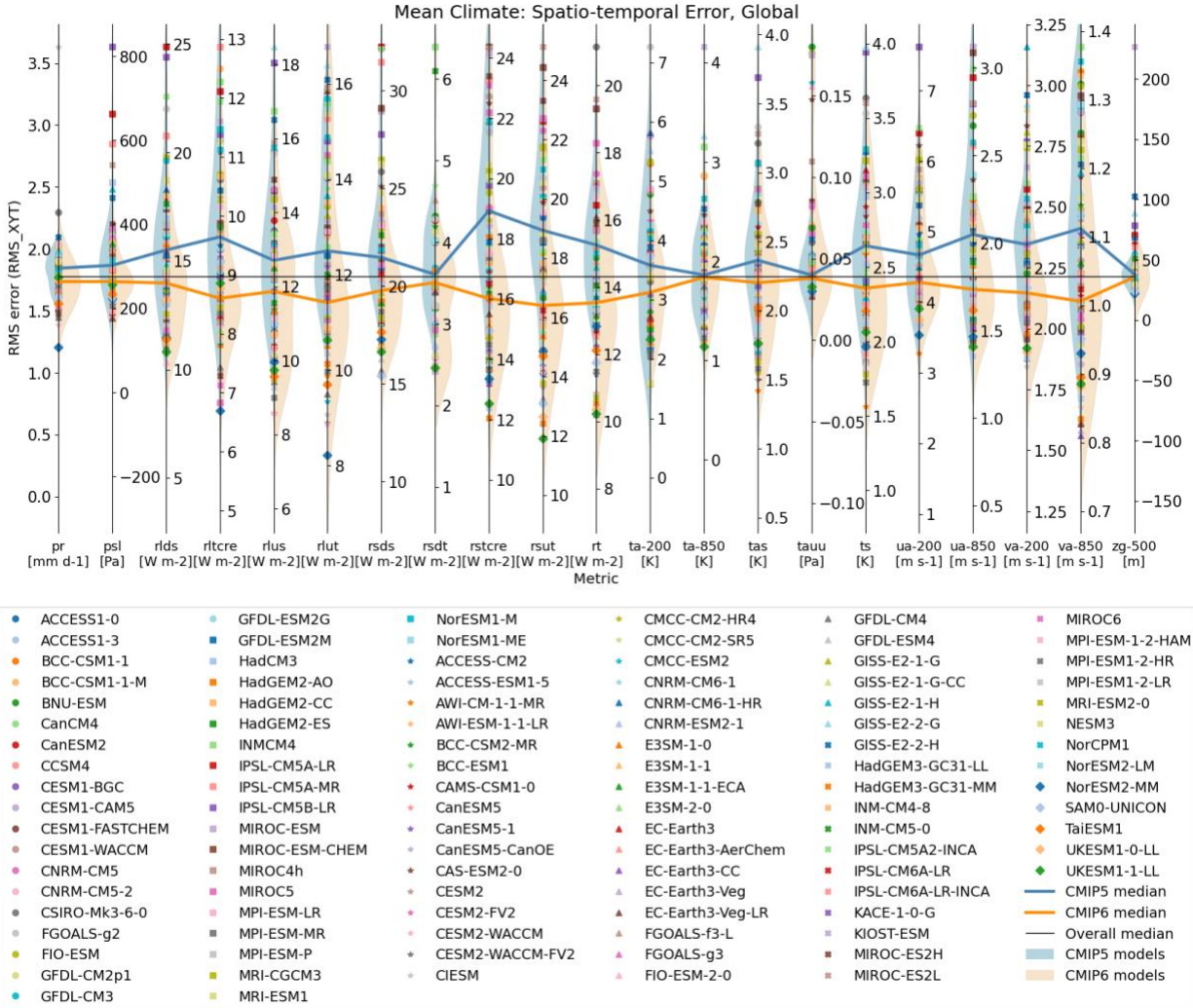
1229 **Table 1.** List of variables and observation datasets used as reference datasets for the PMP's
 1230 mean climate evaluation in this paper (Section 3.1 and Figs. 1-2). A ditto mark (“) indicates the
 1231 same as above.
 1232

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	"	



1233
 1234 **Figure 1.** Portrait plot for spatial RMSE (uncentered) of global seasonal climatologies for (a)
 1235 CMIP5 (models ACCESS1-0 to NorESM1-ME on the ordinate) and (b) CMIP6 (models
 1236 ACCESS-CM2 to UKESM1-1-LL on the ordinate) for 1981-2005 epoch. The RMSE is calculated
 1237 for each season (shown as triangles in each box) over the globe including both land and ocean,
 1238 and model and reference data were interpolated to a common 2.5x2.5 degree grid. The RMSE
 1239 of each variable is normalized by the median RMSE of all CMIP5 and 6 models. A result of 0.2
 1240 (-0.2) is indicative of an error that is 20% greater (lesser) than the median RMSE across all
 1241 models. Models in each group are sorted in alphabetical order. Full names of variable names on
 1242 the abscissa and their reference datasets can be found in Table 1. Detailed information for

1243 models can be found at the *Earth System Documentation* (ES-DOC, <https://search.es-doc.org/>;
1244 Pascoe et al., 2020). The interactive version of the Portrait plot in this figure is available on the
1245 PMP result pages on the PCMDI website (https://pcmdi.llnl.gov/metrics/mean_clim/).



1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258

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/).



1259
 1260 **Figure 3.** Application of ENSO metrics to CMIP6 models. Model names with an asterisk (*)
 1261 indicate that 10 or more ensemble members were used in this analysis. Dots indicate metric
 1262 values from individual ensemble members while bars indicate the average of metric values
 1263 across the ensemble members. Bars colored for easier identification of model names at the
 1264 bottom of the figure. Metrics were grouped into three *Metric Collections*: (a-n) ENSO
 1265 Performance, (o-r) ENSO Teleconnections, and (s-w) ENSO processes. Names of individual
 1266 metrics and default reference datasets being used are noted on top of each panel, and
 1267 observational uncertainty by applying the metrics for alternative reference datasets noted on the
 1268 upper right of each panel is shown as gray-shaded. Detailed descriptions for each metric can be
 1269 found at https://github.com/CLIVAR-PRP/ENSO_metrics/wiki.

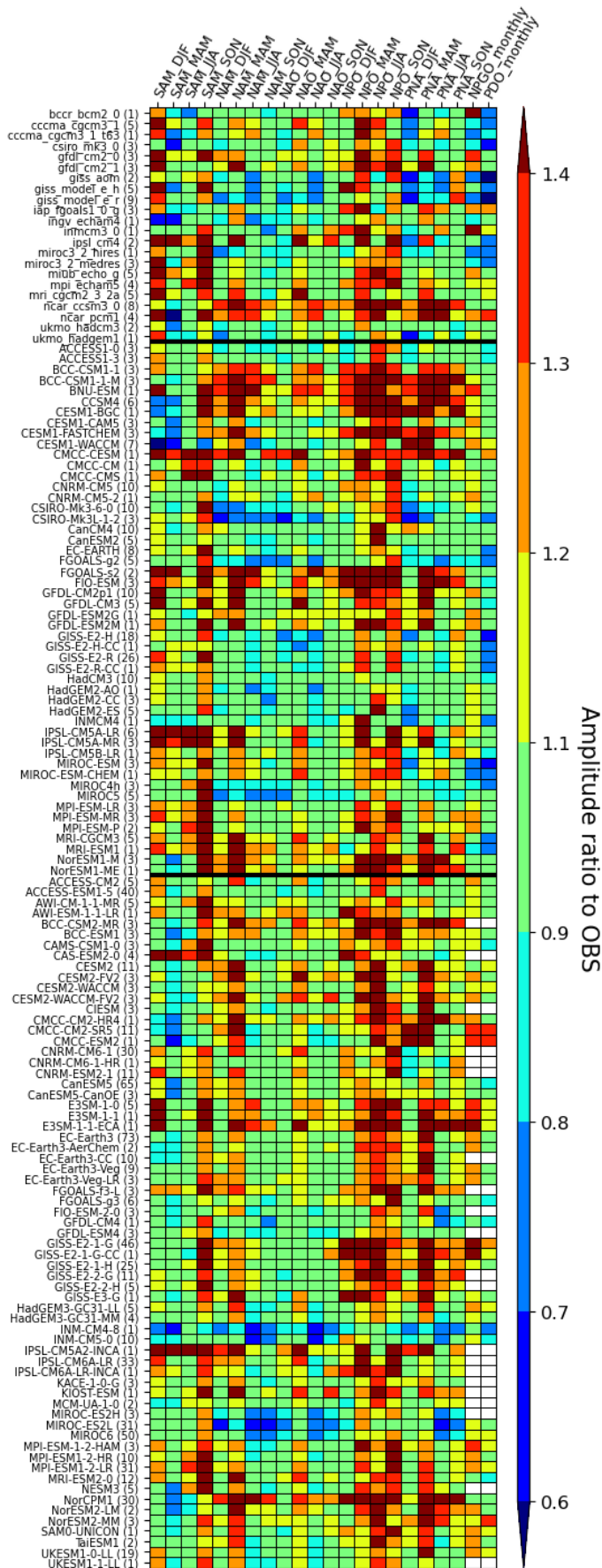
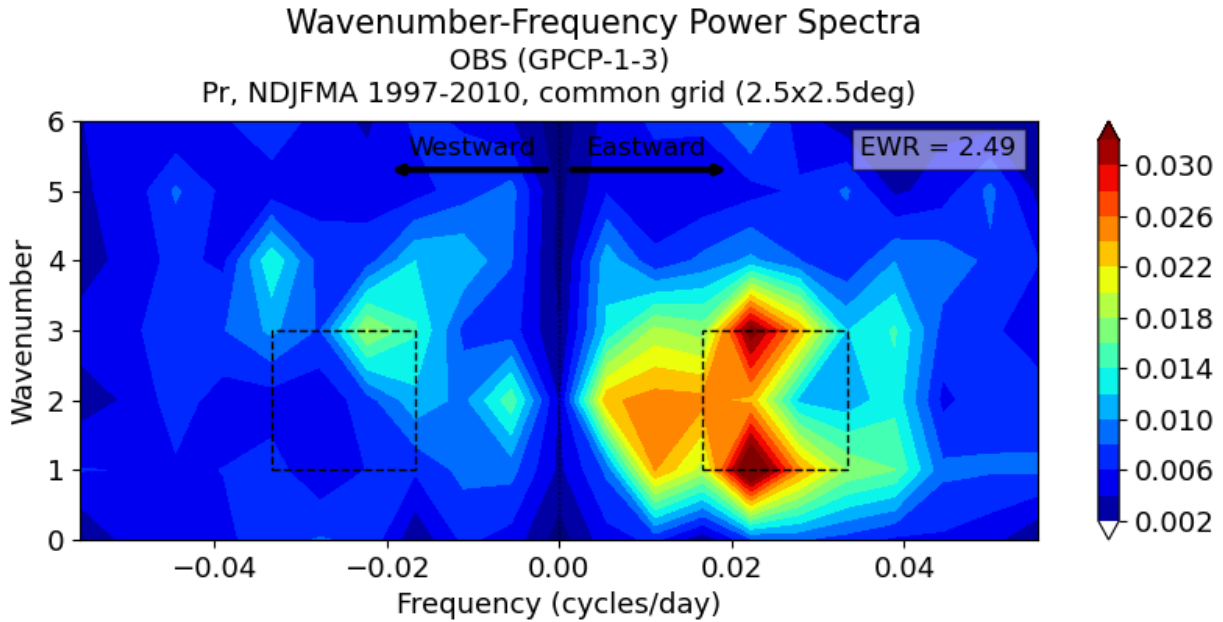
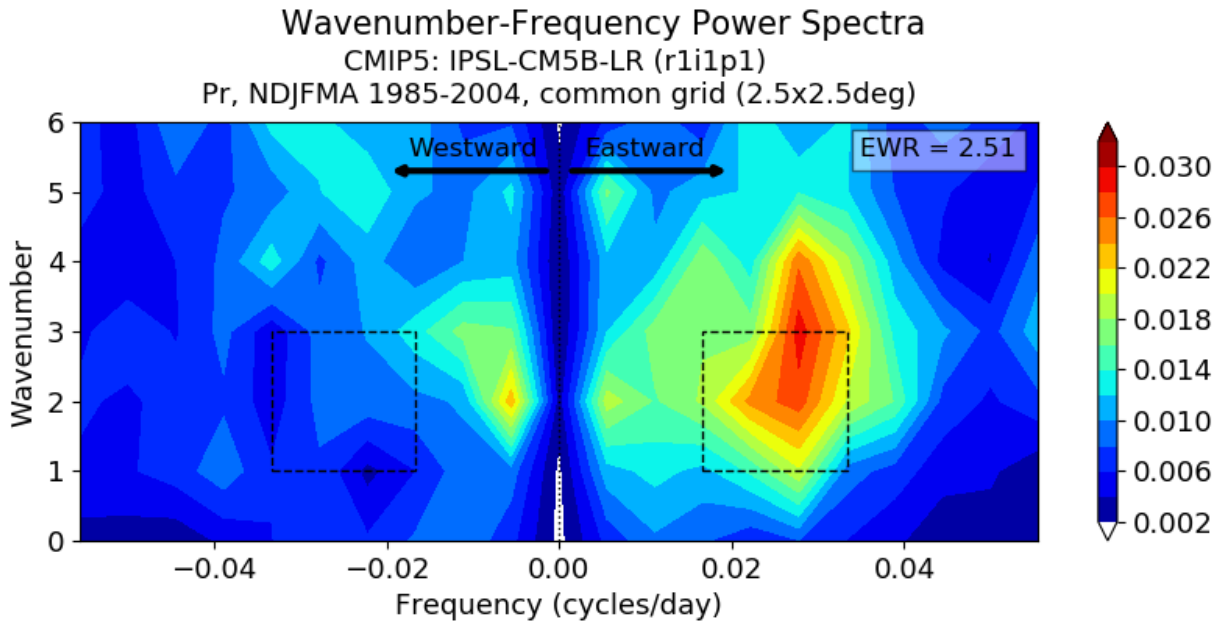


Figure 4. Portrait plots of the amplitude of extratropical modes of variability simulated by CMIP3, 5, and 6 models in their historical or equivalent simulations, as gauged by the ratio of spatiotemporal standard deviations of the model and observed PCs, obtained using the CBF method in the PMP. Columns (horizontal axis) are for mode and season, and rows (vertical axis) are for models from CMIP3 (top), CMIP5 (middle), and CMIP6 (bottom), separated by thick black horizontal lines. For sea level pressure–based modes (SAM, NAM, NAO, NPO, and PNA) in the upper-left hand triangle the model results are shown relative to NOAA-20CR. For SST-based modes (NPGO and PDO), results are shown relative to HadISSTv1.1. Numbers in parentheses following model names indicate the number of ensemble members for the model. Metrics for individual ensemble members were averaged for each model. White boxes indicate missing value.

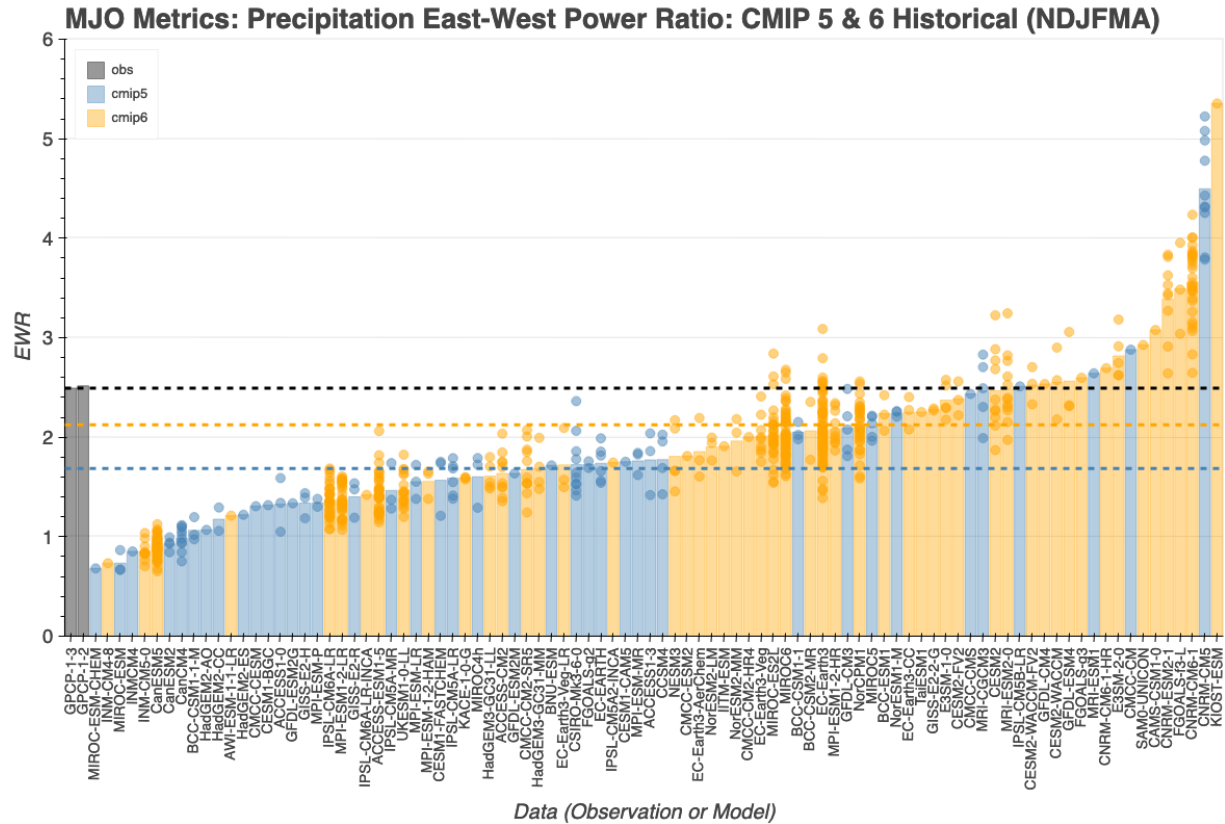
1296
1297 (a) Observation



1298
1299 (b) Model

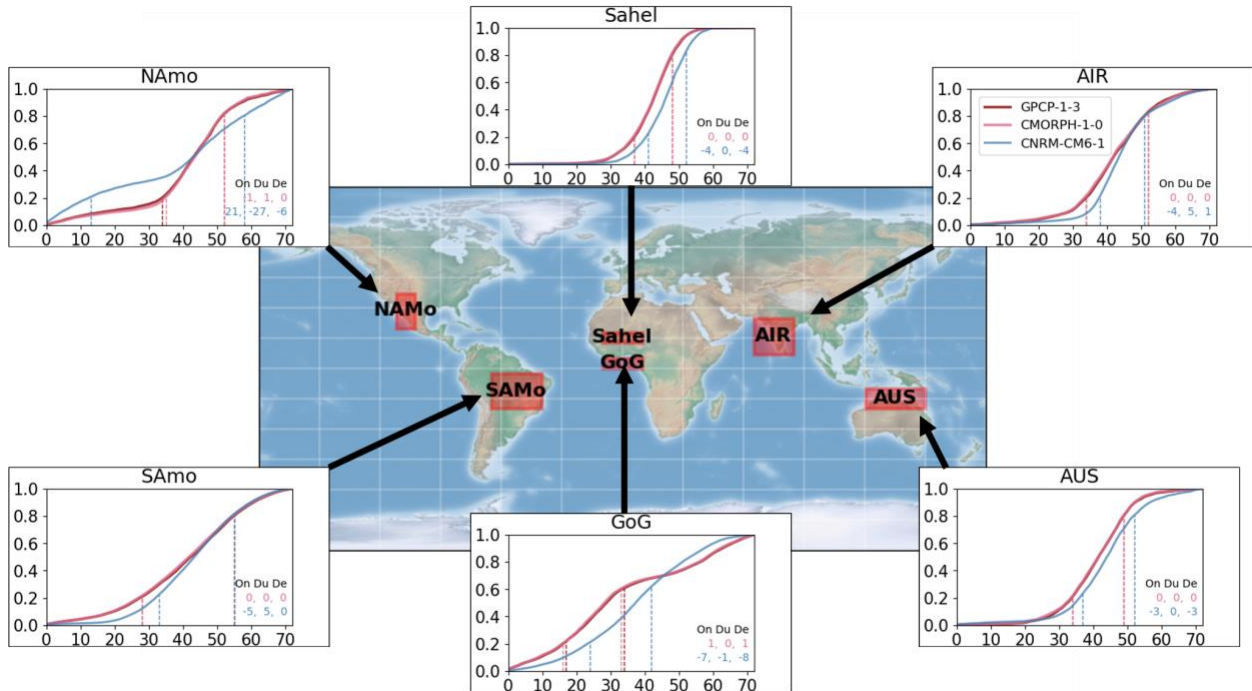


1300
1301
1302 **Figure 5.** MJO EWR diagnostics – wavenumber-frequency power spectra – from (a) GPCP v1.3
1303 (Huffman et al., 2001) and (b) IPSL-CM5B-LR model of CMIP5. The EWR is defined as the ratio
1304 of eastward power (averaged in the box on the right) to westward power (averaged in the box
1305 on the left) from the 2-dimensional wavenumber-frequency power spectra of daily 10°S–10°N
1306 averaged precipitation in November to April (shaded, $\text{mm}^2 \text{day}^{-2}$). Power spectra are calculated
1307 for each year and then averaged over all years of data. The units of power spectra for the
1308 precipitation is $\text{mm}^2 \text{day}^{-2}$ per frequency interval per wavenumber interval.



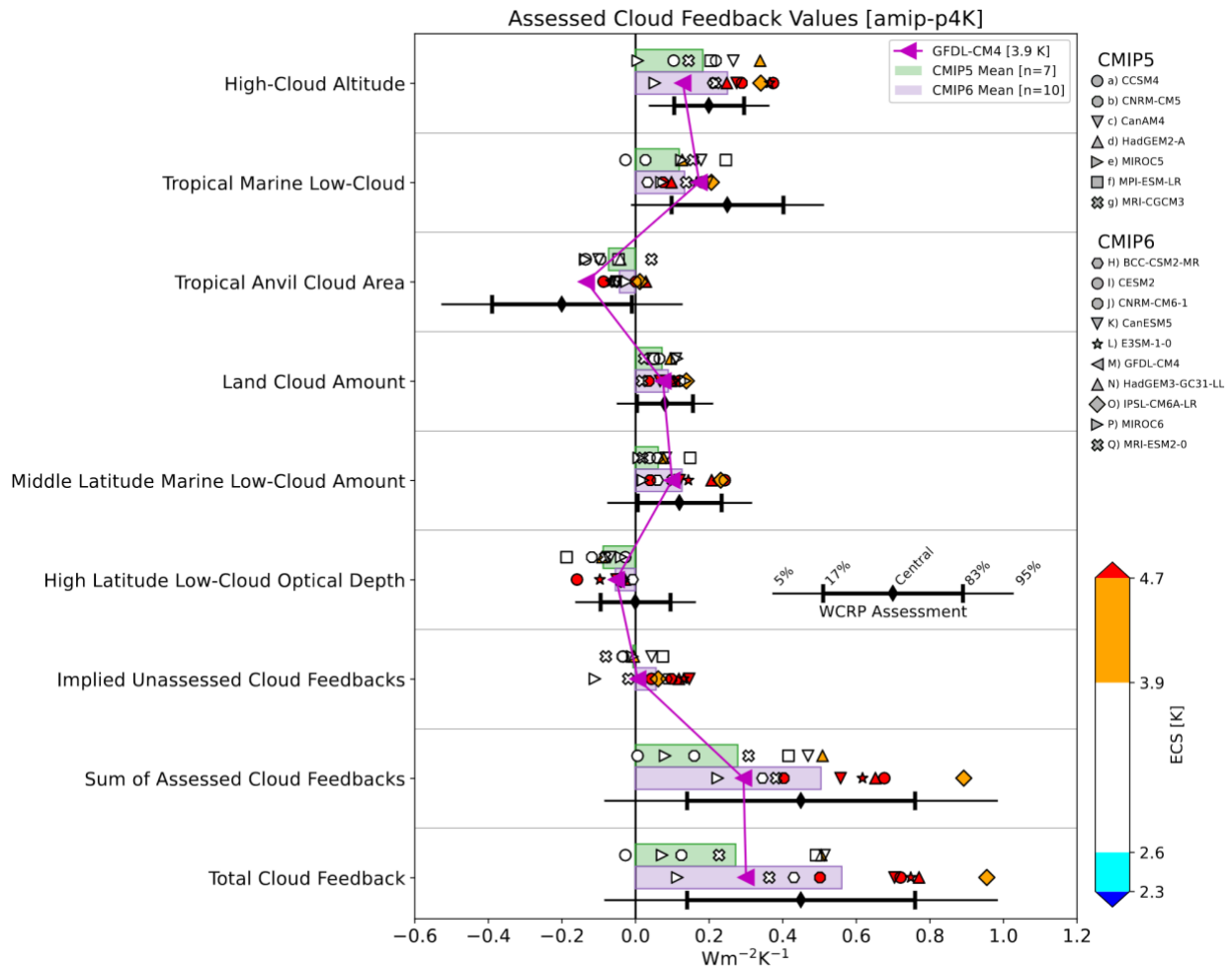
1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319

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.



1320
 1321 **Figure 7.** Demonstration of the monsoon metrics obtained from observation datasets (GPCP
 1322 v1.3 and CMORPH v1.0 (Joyce et al., 2004; Xie et al., 2017)) and a CMIP6 model's Historical
 1323 simulation conducted using CNRM-CM6-1. The results are obtained for monsoon regions: All-
 1324 India Rainfall (AIR), Sahel, Gulf of Guinea (GoG), North American Monsoon (NAM), South
 1325 American Monsoon (SAM), and Northern Australia (AUS). The regions are defined in Sperber
 1326 and Annamalai (2014). Metrics for onset (On), Duration (Du), and Decay (De) derived as
 1327 differences to the default observation (GPCP v1.3) in pentad indices (observation minus model)
 1328 are shown at lower right of each panel. Pentad indices for onset and decay of each region are
 1329 also shown as vertical lines.

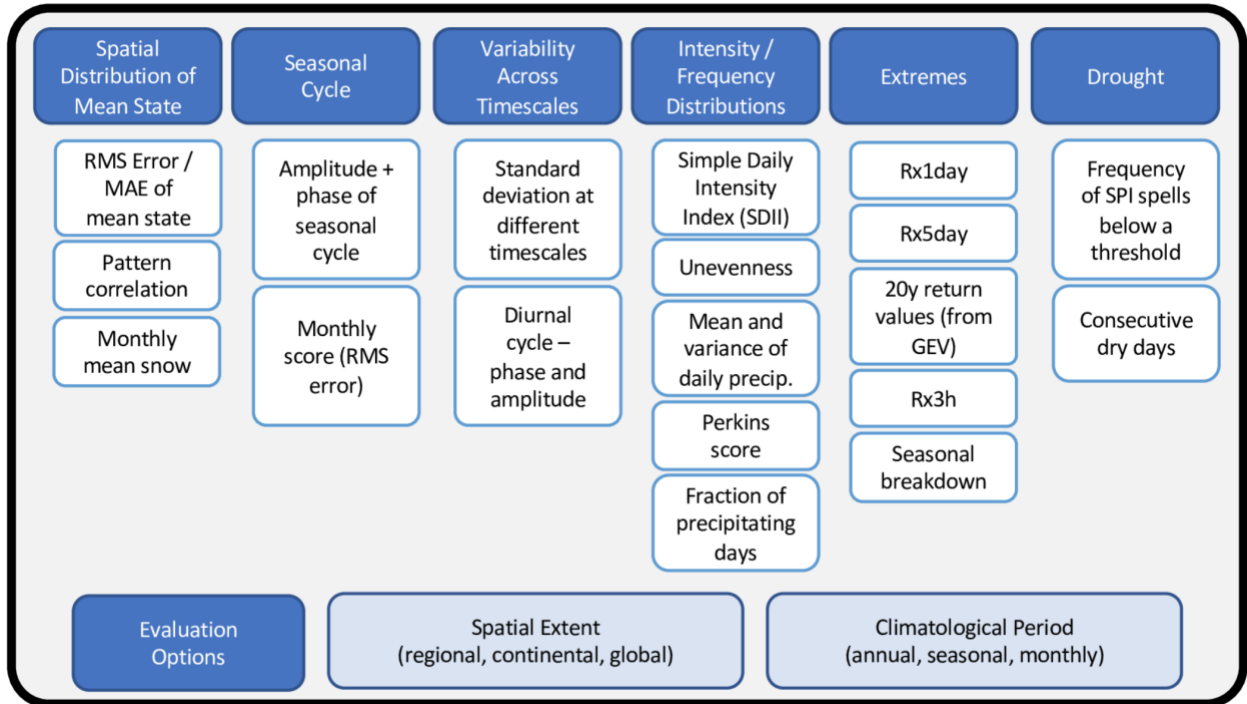
1330



1331
1332
1333
1334
1335
1336
1337

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

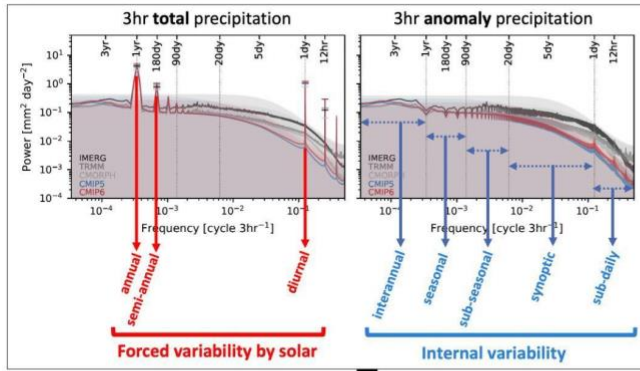
1338



1339
1340
1341
1342
1343
1344
1345

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

(a) Power spectra (Tropics)

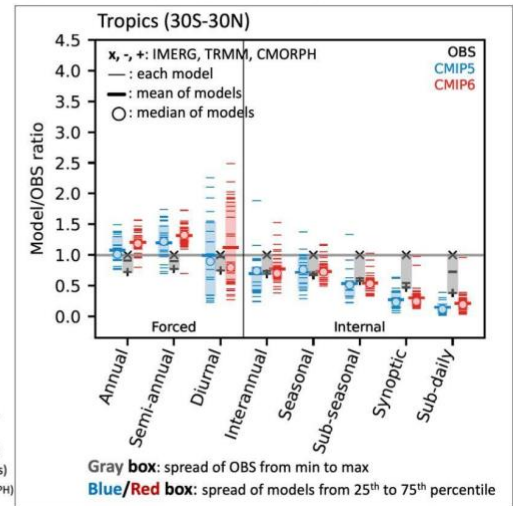


$$\text{Metric} = P_{\text{MODEL}}/P_{\text{OBS}}$$

P: selected or band-averaged power

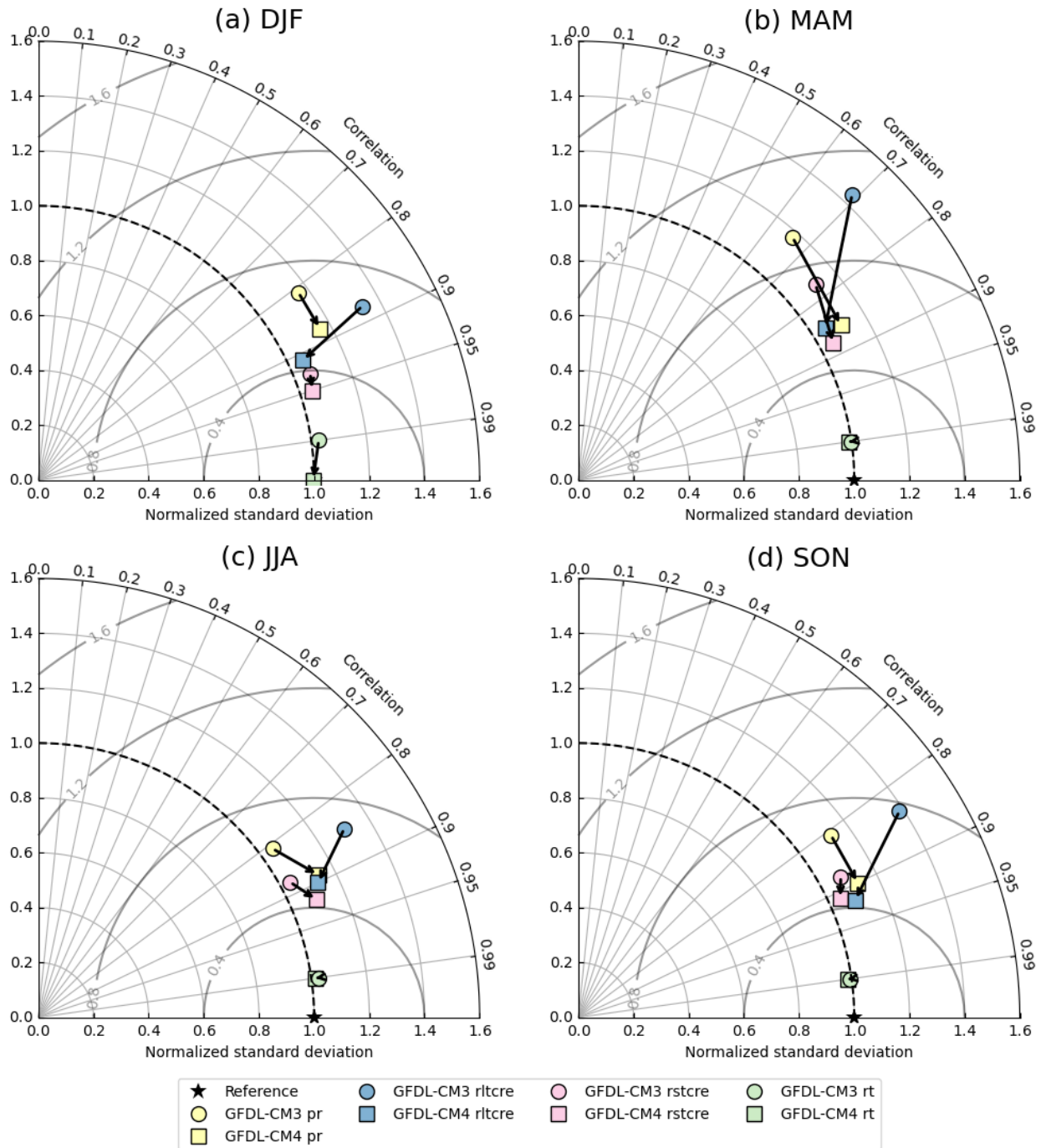
21 CMIP5 (53 realizations)
33 CMIP6 (143 realizations)
3 OBS (IMERG, TRMM, CMORPH)

(b) Metric for precip variability across timescales



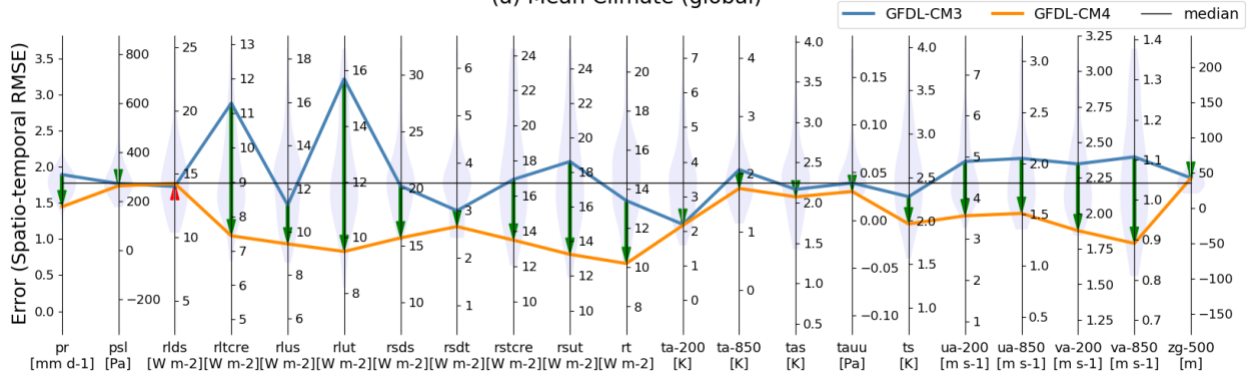
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360

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).



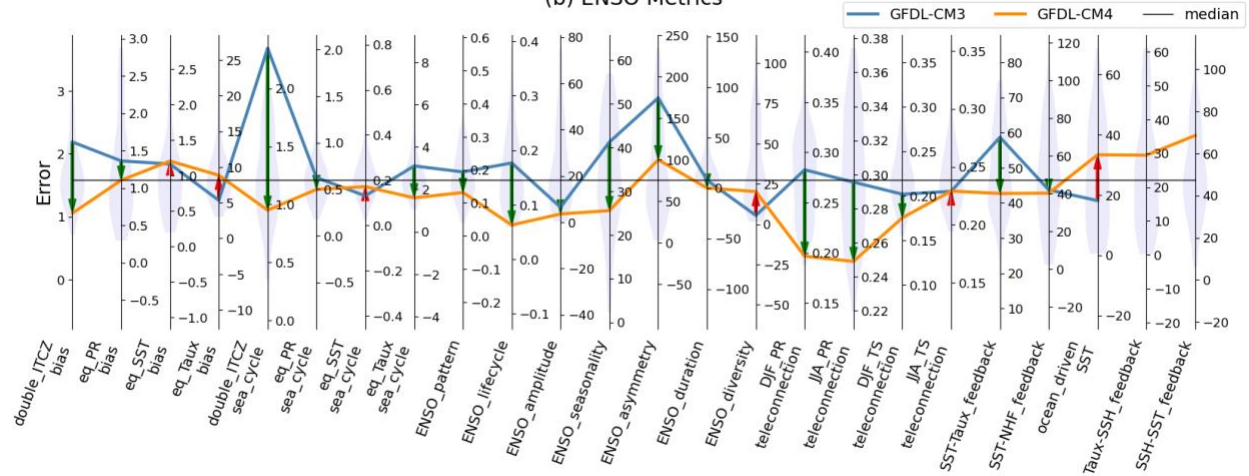
1361
 1362 **Figure 11.** Taylor Diagram contrasting performance of an ESM in their two different versions
 1363 (i.e., GFDL-CM3 from CMIP5 and GFDL-CM4 from CMIP6) in its Historical simulation for
 1364 multiple variables (precipitation [pr], longwave cloud radiative effect [rltcre], shortwave cloud
 1365 radiative effect [rstcre], and total radiation flux [rt]) in their climatology over the globe for (a) DJF,
 1366 (b) MAM, (c) JJA and (d) SON seasons. The arrow is directed toward the newer version of the
 1367 model from the older version (i.e., GFDL-CM3 → GFDL-CM4).

(a) Mean Climate (global)



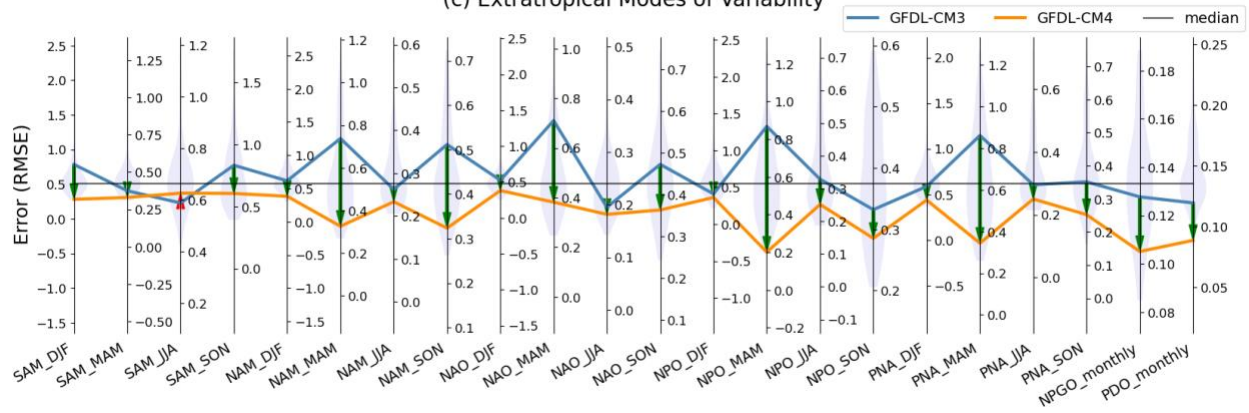
1368

(b) ENSO Metrics



1369

(c) Extratropical Modes of Variability



1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

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.