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

3

4  Jiwoo Lee[1], Peter J. Gleckler[1], Min-Seop Ahn[2,3], Ana Ordonez[1], Paul A. Ullrich[1,4], Kenneth R.

5  Sperber[1,a], Karl E. Taylor[1], Yann Y. Planton[5,6], Eric Guilyardi[7,8], Paul Durack[1], Celine Bonfils[1],

6  Mark D. Zelinka[1], Li-Wei Chao[1], Bo Dong[1], Charles Doutriaux[1], Chengzhu Zhang[1], Tom Vo[1],

7  Jason Boutte[1], Michael F. Wehner[9], Angeline G. Pendergrass[10,11], Daehyun Kim[12], Zeyu Xue[13],

8  Andrew T. Wittenberg[14], and John Krasting[14]

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

27

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 PCMDI Metrics Package (PMP) is an open-source Python software package that provides

38 "quick-look" objective comparisons of ESMs with one another and with observations. The comparisons include

39 metrics of large- to global-scale climatologies, tropical inter-annual and intra-seasonal variability modes such as El

40 Niño-Southern Oscillation (ENSO) and Madden-Julian Oscillation (MJO), extratropical modes of variability, regional

41 monsoons, cloud radiative feedbacks, and high-frequency characteristics of simulated precipitation, including

42 extremes. The PMP results are produced in the context of all model simulations contributed to CMIP6 and earlier

43 CMIP phases. An important priority of the PMP is to document evaluation statistics for all Historical and AMIP

44 simulations submitted to recent phases of CMIP, providing version-controlled information for all data sets and

45 software packages being used. Among other purposes, this also enables modeling groups to assess performance

46 changes during the ESM development cycle in the context of the error distribution of the multi-model ensemble. In

47 this paper, we present an overview of the PMP including its history to date, capabilities, recent updates, and future

48 direction.

## 1 Introduction

Earth System Models (ESMs) are key tools for projecting climate change and conducting research to enhance our understanding of the Earth system. Enhancing the reliability of models is therefore important, yet evaluating ESMs is a complex endeavor, given the vast range of climate characteristics across space and time scales. A necessary step to evaluate the performance of ESMs is quantifying their consistency with available observations.

The Program for Climate Model Diagnosis and Intercomparison (PCMDI) has worked closely with the World Climate Research Programme's (WCRP) Working Group on Coupled Models (WGCM) and Working Group on Numerical Experimentation (WGNE) to design and support Model Intercomparison Projects (MIPs) (Potter et al., 2011). This effort began with the Atmospheric Model Intercomparison Project (AMIP; Gates, 1992; Gates et al., 1999), and has continued through multiple phases of the Coupled Model Intercomparison Project (CMIP; Meehl et al., 1997, 2000, 2007; Covey et al., 2003; Taylor et al., 2012). The most recent phase of CMIP (CMIP6; Eyring et al., 2016) provides a set of well-defined experiments that most climate modeling centers perform, and subsequently makes results available for a large and diverse community to analyze.

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

Considering the exponential growth of data size and diversity of ESM simulations, there has been a pressing need for the research community to become more efficient and systematic in evaluating ESMs and documenting their performances. To respond to the need, PCMDI has developed the PCMDI Metrics Package (PMP), to quantitatively synthesize results from the archive of CMIP simulations via performance metrics that help characterize the overall agreement between models and observations (Gleckler et al., 2016). In this paper, we describe the latest update of the PMP and its focus on providing a diverse suite of summary statistics that can be used to construct "quick-look" summaries of ESM performance from simulations made publicly available to the research community, notably CMIP. For our purposes, "performance metrics" are typically (but not exclusively) well-established statistical measures that

86    quantify the consistency between observed and simulated characteristics. One goal of the PMP is to further diversify

87    the suite of high-level performance tests that help characterize the simulated climate. The results provided by the PMP

88    are frequently used to address two overarching and recurring questions: 1) What are the relative strengths and

89    weaknesses between different models? and 2) How are models improving with further development? Addressing the

90    second question is often referred to as "benchmarking" and this motivates an important emphasis of the effort

91    described in this paper—striving to advance the documentation of all data and results of the PMP in an open and

92    ultimately reproducible manner.

93        The rest of the paper is organized as follows. In section 2, we provide a technical description of the PMP and

94    its accompanying reference datasets. In section 3, we describe various sets of simulation metrics that capture an

95    increasingly comprehensive range of physical processes and time scales ranging from hours to centurial. In section 4,

96    we introduce the usage of PMP for model benchmarking. In section 5, we discuss the remaining challenges, and we

97    conclude in section 6 with a summary and future direction.

98

99    **2 Software package and data description**

100   The PMP is a Python-based open-source software framework (https://github.com/PCMDI/pcmdi_metrics) designed

101   to objectively gauge the consistency between ESMs and available observations via well-established statistics. The

102   PMP has been mainly used for the evaluation of CMIP-class models. A subset of CMIP experiments are particularly

103   well suited to comparing models with observations. The experiments of particular interest include those involving

104   prescribed sea surface temperature (SST) in accordance with the AMIP protocol, as well as coupled model simulations

105   labeled as "Historical" that are driven by varying natural and anthropogenic forcings. Some of the metrics applicable

106   to these experiments may also be relevant to others (e.g., multi-century coupled control runs called "PiControl" and

107   idealized "4xCO2" simulations that are designed for estimating climate sensitivity).

108       The PMP has been applied to multiple generations of CMIP in a quasi-operational fashion as new simulations

109   are made available, new analysis methods are incorporated, or new observational data become accessible. Shortly

110   after simulations from the most recent phase of the CMIP (i.e., CMIP6) became accessible, PMP quick-look

111   summaries were provided on the PCMDI's website (https://pcmdi.llnl.gov/metrics/), offering a resource to scientists

112   involved in CMIP or others interested in the evaluation of ESMs. To facilitate this, in PCMDI the PMP is technically

113   linked to the Earth System Grid Federation (ESGF) that is a primary CMIP data delivery infrastructure (Williams et

114   al., 2016).

115       The PMP is designed to readily work with model output that has been processed using the Climate Model

116   Output Rewriter (CMOR; https://cmor.llnl.gov/), which is a software library developed to prepare model output as

117   CF-compliant (Hassell et al., 2017; Eaton et al., 2022, http://cfconventions.org/) netCDF files. The CMOR is used by

118   most modeling groups contributing to CMIP, ensuring all model output adheres to the CMIP data structures that

119   themselves are based on the CF conventions. It is possible to use the PMP on model output that has not been prepared

120   by CMOR, but this usually requires additional work, e.g., mapping the data to meet the community standards.

121       For reference datasets, the PMP uses observational products processed to be compliant with the Observations

122   for Model Intercomparison Projects (obs4MIPs; https://pcmdi.github.io/obs4MIPs/). The obs4MIPs effort was

123 initiated circa 2010 (Gleckler et al., 2011) to advance the use of the observations in model evaluation and research.

124 Substantial progress has been made in establishing obs4MIPs data standards that technically align with CMIP model

125 output (e.g., Teixeira et al., 2014; Ferraro et al., 2015), with the data products published on the ESGF (Waliser et al.,

126 2020). Obs4MIPs-compliant data were prepared with CMOR, and the data directly available via obs4MIPs are used

127 as PMP reference datasets.

128 The PMP leverages other Python-based open-source libraries. A primary fundamental tool used in the latest

129 PMP version is the Python package, Xarray Climate Data Analysis Tools (xCDAT; Vo et al., 2023;

130 https://xcdat.readthedocs.io). The xCDAT is developed to provide a more efficient, robust, and streamlined user

131 experience in climate data analysis when using xarray (https://docs.xarray.dev/). Portions of the PMP rely on the

132 precursor of the xCDAT, a Python library called Community Data Analysis Tools (CDAT, Williams et al., 2009;

133 Williams, 2014; Doutriaux et al., 2019), which has been fundamental since the early development stages of the PMP.

134 The xarray software provides much of the functionality of CDAT (e.g., I/O, indexing, and subsetting). However, it

135 lacks some key climate domain features that have been frequently used by scientists and exploited by the PMP (e.g.,

136 regridding, utilization of spatial/temporal bounds for computational operations) which motivated the development of

137 the xCDAT. Completing the transition from CDAT to xCDAT is a technical priority for the next version of PMP.

138 The primary delivery output of the PMP is the summary statistics. We strive to make the baseline results (raw

139 statistics) publicly available and well-documented, and continue to make advances with this priority. For our purposes,

140 we are referring to model performance "summary statistics" and "metrics" interchangeably, although in some

141 situations we consider there to be an important distinction. For us, a genuine performance metric constitutes a well-

142 defined and established statistic that has been used in a very specific way (e.g., a particular variable, analysis, and

143 domain) for long-term benchmarking (see Section 4). The distinction between summary statistics and metrics is

144 application-dependent and evolving as the community advances efforts to establish quasi-operational capabilities to

145 gauge ESM performance. Some visualization capabilities described in Section 3 are made available through the PMP.

146 Users can also further explore the model data comparisons using their preferred visualization methods or incorporate

147 the results into their own studies from the summary statistics from the PMP. Noting the above, the scope of the PMP

148 is fairly targeted. It is not intended to be "all-purpose", e.g. by incorporating the vast range of diagnostics used in

149 model evaluation.

150 To help advance open and reproducible science, the PMP has been maintained with an open-source policy

151 with accompanying metadata for data reproducibility and reusability. The PMP code is distributed and released with

152 version control. Online documentation (http://pcmdi.github.io/pcmdi_metrics/), including user demo Jupyter

153 Notebooks, and a database of pre-calculated PMP statistics for all AMIP and Historical simulations in the CMIP

154 archive are also available online. The archive of these statistics stored as JSON files (Crockford, 2006; Crockford and

155 Morningstar, 2017) includes versioning details for all codes, and dependencies and data that were used for the

156 calculations. These files provide the baseline results of the PMP (See the Code and Data Availability section for

157 details). Advancements in model evaluation along with the number of models and complexity of simulations motivate

158 more systematic documentation of performance summaries. With PMP workflow provenance information being

159 recorded and the model and observational data standards maintained by PCMDI and colleagues, PMP strives to make
160 all its results reproducible.

161

162 **3 Current PMP capabilities**

163 The PMP builds upon model performance tests that have resulted from research at PCMDI and via close
164 collaborations. Contributors have helped expand the PMP beyond its traditional large-scale performance summaries
165 of the mean climate (Gleckler et al., 2008). Various evaluation metrics have been implemented to the PMP for climate
166 variability such as El Niño-Southern Oscillation (ENSO) (Planton et al., 2021; Lee et al., 2021a), extratropical modes
167 of variability (Lee et al., 2019, 2021b), intra-seasonal oscillation (Ahn et al., 2017), monsoons (Sperber and
168 Annamalai, 2014), cloud feedback (Zelinka et al., 2022), and the characteristics of simulated precipitation
169 (Pendergrass et al., 2020; Ahn et al., 2022, 2023) and extremes (Wehner et al., 2020, 2021). This section will provide
170 an overview of each category of the current PMP evaluation metrics with their usage demonstrations.

171

172 *3.1 Climatology*

173 Mean state metrics quantify how well models simulate observed climatological fields at a large scale, gauged by a
174 suite of well-established statistics that have been used in climate research for decades. The focus is on the coupled
175 "Historical" and atmospheric-only AMIP (Gates et al., 1999) simulations which are well-suited for comparison with
176 observations. The PMP extracts seasonally and annually averaged fields of multiple variables from large-scale
177 observationally based datasets and results from model simulations. Different obs4MIPs-compliant reference datasets
178 are used depending on the variable examined. When multiple reference datasets are available, one of them is
179 considered as a "default" while others are identified as "alternatives". The default datasets are typically state-of-the-
180 art products, but in general, we lack definitive measures as to which is the most accurate, so the PMP metrics are
181 routinely calculated with multiple products so that it can be determined what difference the selection of alternative
182 observations makes to judgment made about model fidelity. The suite of mean climate metrics (all area weighted)
183 includes spatial and spatiotemporal root-mean-square error (RMSE), centered spatial RMSE, spatial-mean bias, spatial
184 standard deviation, spatial pattern correlation, and spatial and spatiotemporal mean absolute error (MAE) of the annual
185 or seasonal climatological time-mean (Gleckler et al., 2008). Often, a space-time statistic is used that gauges both the
186 consistency of the observed and simulated climatological pattern as well as its seasonal evolution (see Eq. 1 from
187 Gleckler et al., 2008). By default, results are available for selected large-scale domains, including: "Global", "Northern
188 Hemisphere (NH) Extratropics" (30ºN-90ºN), "Tropics" (30ºS-30ºN), and "Southern Hemisphere (SH) Extratropics"
189 (30ºS-90ºS). For each domain, results can also be computed for the land and ocean, land only, or ocean only. These
190 commonly used domains highlight the application of the PMP mean climate statistics at large to global scales, but we
191 note that PMP allows users to define their own domains of interest, including at regional scales.

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

196    communities (see Section 5). To interpret results across CMIP models for many variables, we routinely construct
197    normalized Portrait Plots or Gleckler Plots (Gleckler et al., 2008) that provide a quick-look examination of the
198    strengths and weaknesses of different models. For example, in Figure 1, the PMP results display quantitative
199    information of simulated seasonal climatologies of various meteorological model variables via a normalized global
200    spatial RMSE (Gleckler et al., 2008). Variants of this plot have been widely used for presenting model evaluation
201    results, for example, in the Intergovernmental Panel on Climate Change (IPCC) Fifth (Flato et al., 2014, Figures 9.7,
202    9.12, and 9.37) and Sixth Assessment Reports (Eyring et al., 2021, Chapter 3, Figure 3.42). Because the error
203    distribution across models is variable dependent, the statistics are often normalized to help reveal differences, in this
204    case via the median RMSE across all models (see Gleckler et al. 2008 for more details). This normalization enables a
205    common color scale to be used for all statistics on the Portrait Plot, highlighting the relative strengths and weaknesses
206    of different models. In this example (Fig. 1), an error of -0.5 indicates that a model's error is 50% smaller than the
207    typical (median) error across all models, whereas an error of 0.5 is 50% larger than the typical error in the multi-model
208    ensemble. In many cases, the horizontal bands in the Gleckler plots show that simulations from a given modeling
209    center have similar error structures relative to the multi-model ensemble.

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

230

### 3.2 El Niño-Southern Oscillation

231
232 The El Niño-Southern Oscillation (ENSO) is Earth's dominant interannual mode of climate variability, which impacts
233 global climate via both regional oceanic effects and far-reaching atmospheric teleconnections (McPhaden et al., 2006,
234 2020). In response to increasing interest in a community approach to ENSO evaluation in models (Bellenger et al.,
235 2014), the International Climate and Ocean Variability, Predictability and Change (CLIVAR) Research Focus on
236 ENSO in a Changing Climate, together with the CLIVAR Pacific Region Panel, developed the CLIVAR ENSO
237 Metrics Package (Planton et al., 2021) which is now utilized within the PMP. The ENSO metrics are divided into three
238 Metrics Collections: *Performance* (i.e., background climatology and basic ENSO characteristics), *Teleconnections*
239 (ENSO's worldwide teleconnections), and *Processes* (ENSO's internal processes and feedback). Planton et al. (2021)
240 found that CMIP6 models generally outperform CMIP5 models in several ENSO metrics in particular for those related
241 to tropical Pacific seasonal cycles and ENSO teleconnections. This effort is discussed in more detail in Planton et al.
242 (2021), and detailed descriptions of each metric in the package are available in the ENSO Package online open-source
243 code repository on its GitHub Wiki pages (see https://github.com/CLIVAR-PRP/ENSO_metrics/wiki).

244       Figure 3 demonstrates the application of the ENSO metrics to CMIP6, showing the magnitudes of inter-
245 model and inter-ensemble spreads, along with observational uncertainty varying across metrics. For a majority of the
246 ENSO Performance metrics model error and inter-model spread are substantially larger than observational uncertainty
247 (Figs. 3a-n). This highlights the systematic biases like the double ITCZ (Fig. 3a) that are persisting through CMIP
248 phases (Tian and Dong, 2020). Similarly, ENSO Processes metrics (Figs. 3t-w) indicate large errors in the feedback
249 loops generating SST anomalies, indicating a different balance of processes in the model and in the reference and
250 possibly compensating errors (Bayr et al., 2019, Guilyardi et al. 2020). In contrast, for ENSO Teleconnection metrics,
251 the observational uncertainty is substantially larger, thus challenging validation of model error (Figs. 3o-r). For some
252 metrics, such as the ENSO duration (Fig. 3f), the ENSO Asymmetry metric (Fig. 3i), and the Ocean driven SST metric
253 (Fig. 3s), there are larger inter-ensemble spreads than the inter-model spreads. From such results, Lee et al. (2021a)
254 examined the inter-model and inter-member spread of these metrics from the large ensembles available from CMIP6
255 and the US CLIVAR Large Ensemble Working Group. They argued that to robustly characterize baseline ENSO
256 characteristics and physical processes, larger ensemble sizes are needed, compared to existing state-of-the-art
257 ensemble projects.

258

### 3.3 Extratropical Modes of Variability

259
260 The PMP includes objective measures of the pattern and amplitude of extratropical modes of variability from
261 PCMDI's research, which has expanded beyond its traditional large-scale performance summaries to include
262 interannual variability, considering increasing interest in setting an objective approach for the collective evaluation of
263 multiple modes. Extratropical modes of variability (ETMoV) metrics in the PMP were developed by Lee et al. (2019a)
264 that stem from earlier works (e.g., Stoner et al., 2009; Phillips et al., 2014). Lee et al. (2019a) illustrated a challenge
265 when evaluating modes of variability using the traditional empirical orthogonal functions (EOF). In particular, when
266 a higher-order EOF of a model more closely corresponds to a lower-order observationally based EOF (or vice versa),
267 it can significantly affect conclusions drawn about model performance. To circumvent this issue in evaluating the

268     interannual variability modes, Lee et al. (2019a) used the Common Basis Function (CBF) approach that projects the

269     observed EOF pattern onto model anomalies. This approach has been previously applied for the evaluation of

270     intraseasonal variability modes (Sperber, 2004; Sperber et al., 2005), and recently for Antarctic climate change (Jun

271     et al., 2020), seasonal-to-decadal predictability associated with the ENSO (Choi and Son, 2022). In the PMP, the CBF

272     approach is taken as a default method, and the traditional EOF approach is also enabled as an option for the ETMoV

273     metrics calculations.

274     The ETMoV metrics in the PMP measure simulated patterns and amplitudes of ETMoV, and quantify their

275     agreement with observations (e.g., Lee et al., 2019a, 2021b). The PMP's ETMoV metrics evaluate 5 atmospheric

276     modes – the Northern Annular Mode (NAM), North Atlantic Oscillation (NAO), Pacific North America pattern

277     (PNA), North Pacific Oscillation (NPO), and Southern Annular Mode (SAM), and 3 ocean modes diagnosed by the

278     variance of sea-surface temperature – Pacific Decadal Oscillation (PDO), North Pacific Gyre Oscillation (NPGO),

279     and Atlantic Multi-decadal Oscillation (AMO). The AMO is included for experimental purposes, considering the

280     significant uncertainty in detecting the AMO (Deser and Philips 2021; Zhao et al., 2022). The amplitude metric,

281     defined as the ratio of standard deviations of the model and observed principal components, has been used to examine

282     the evolution of the performance of models across different CMIP generations (Fig. 4, adapted from Lee et al., 2021b).

283     Green shading predominates, indicating where the simulated amplitude of variability is similar to observations. In

284     some cases, such as for SAM_SON, the models overestimate the observed amplitude. Other authors have used Portrait

285     plots to synthesize CMIP performance of simulated variability (e.g., Sillmann et al., 2013; Bellenger et al., 2014;

286     Cannon 2020; Kim et al., 2020; Planton et al., 2020; Zhang et al., 2021; Ahn et al., 2022, 2023).

287     The PMP's ETMoV metrics have been used in several model evaluation studies. For example, Orbe et al.

288     (2020) analyzed models from U.S. climate modeling groups including DOE, National Aeronautics and Space

289     Administration (NASA), National Center for Atmospheric Research (NCAR), and National Oceanic and Atmospheric

290     Administration (NOAA), where they found that the improvement in the ETMoV performance is highly dependent on

291     mode and season, when comparing across different generations of those models. Sung et al. (2021) examined the

292     performance of models run at the Korea Meteorological Administration (K-ACE and UKESM1) in reproducing

293     ETMoVs from their Historical simulations, and concluded that these models reasonably capture most ETMoVs. Lee

294     et al. (2021b) collectively evaluated ~130 models from CMIP3, 5, and 6 archive databases using their ~850 Historical

295     and ~300 AMIP simulations, where they found the spatial pattern skill improved in CMIP6 compared to CMIP5 or

296     CMIP3 for most modes and seasons, while the improvement in amplitude skill is not clear. Arcodia et al. (2023) used

297     the PMP to derive PDO and AMO to investigate their role in decadal variability of subseasonal predictability of

298     precipitation over the western coast of North America and concluded that no significant relationship was found.

299

### 3.4 Intraseasonal Oscillation

301 The PMP has implemented metrics for the Madden-Julian Oscillation (MJO; Madden and Julian, 1971, 1972, 1994).

302 The MJO is the dominant mode of tropical intraseasonal variability, characterized by a pronounced eastward

303 propagation of large-scale atmospheric circulation coupled with convection with a typical periodicity of 30-60 days.

304    Selected metrics from the MJO diagnostics package, developed by the CLIVAR MJO Working Group (Waliser et al.,

305    2009), have been implemented in the PMP following Ahn et al. (2017).

306        We particularly focused on a metric called East/West power Ratio (hereafter, EWR) and East power

307    normalized by Observation (hereafter, EOR). The EWR, proposed by Zhang and Hendon (1997), is defined as the

308    ratio of the total spectral power over the MJO band (eastward propagating, wavenumber 1-3 and period of 30-60 days)

309    to that of its westward propagating counterpart in the wavenumber-frequency power spectra. The EWR metric has

310    been widely used in the community, to examine the robustness of the eastward propagating feature of the MJO (e.g.,

311    Zhang and Hendon, 1997; Hendon et al., 1999; Lin et al., 2006; Kim et al., 2009; Ahn et al., 2017). The EOR is

312    formulated by normalizing a model's spectral power within the MJO band by the corresponding observed value. Ahn

313    et al. (2017) showed EWRs and EORs of the CMIP5 models. Using daily precipitation, the PMP calculates EWR and

314    EOR separately for boreal winter (November to April) and boreal summer (March to October). We apply the

315    frequency-wavenumber decomposition method to precipitation from observations (GPCP-based; 1997-2010) and the

316    CMIP5 and CMIP6 Historical simulations for 1985-2004. For disturbances with wavenumbers 1-3 and frequencies

317    corresponding to 30-60 days, it is clear in observations that the eastward propagating signal dominates over its

318    westward propagating counterpart with an EWR value of about 2.49 (Fig. 5a). Figure 5b shows the wavenumber-

319    frequency power spectrum from CMIP5 IPSL-CM5B-LR as an example, which has an EWR value that is comparable

320    to the observed value.

321        Figure 6 shows the EWR from individual models' multiple ensemble members and their average. The average

322    EWR of the CMIP6 model simulations is more realistic than that of the CMIP5 models. Interestingly, a substantial

323    spread exists across models and also among ensemble members of a single model. For example, while the average

324    EWR value for the CESM2 ensemble is 2.47 (close to 2.49 from GPCP observations), the EWR values of the

325    individual ensemble members range from 1.87 to 3.23. Kang et al. (2020) suggested that the ensemble spread in the

326    propagation characteristics of the MJO can be attributed to the differences in the moisture mean state, especially its

327    meridional moisture gradient. A cautionary note should be given to the fact that the MJO frequency and wavenumber

328    windows are chosen to capture the spectral peak in observations. Thus, while the EWR provides an initial evaluation

329    of the propagation characteristics of the observed and simulated MJO, it is instructive to look at the frequency-

330    wavenumber spectra, as in some cases the dominant periodicity and wavenumber in a model may be different than in

331    observations. It is worthwhile to note that the PMP can be used to obtain EWR and EOR of other daily variables for

332    MJO analysis, such as outgoing longwave radiation (OLR) or zonal wind at 850 hPa (U-850) or 250 hPa (U-250), as

333    shown in Ahn et al. (2017).

334

335    ***3.5 Monsoons***

336    Based on the work of Sperber and Annamalai (2014), skill metrics in the PMP quantify how well models represent

337    the onset, decay, and duration of regional monsoons. From observations and Historical simulations, the climatological

338    pentads of precipitation are area-averaged for six monsoon-related domains: All-India Rainfall, Sahel, Gulf of Guinea,

339    North American Monsoon, South American Monsoon, and Northern Australia, as seen in Fig. 7. For the domains in

340    the Northern Hemisphere, the 73 climatological pentads run from January to December, while for the domains in the

341    Southern Hemisphere, the pentads run from July to June. For each domain, the precipitation is accumulated at each

342    subsequent pentad and then divided by the total precipitation to give the fractional accumulation of precipitation as a

343    function of pentad. Thus, the annual cycle behavior is evaluated irrespective of whether a model has a dry or wet bias.

344    Except for GoG, the onset and decay of monsoon occur for a fractional accumulation of 0.2 and 0.8, respectively.

345    Between these fractional accumulations, the accumulation of precipitation is nearly linear as the monsoon season

346    progresses. Comparison of the simulated and observed onset, duration, and decay are presented in terms of the

347    difference in the pentad index obtained from the model and observations (i.e., model minus observations). Therefore,

348    negative values indicate that the onset or decay in the model occurs earlier than in observations, while positive values

349    indicate the opposite. For duration, negative values indicate that for the model it takes fewer pentads to progress from

350    onset to decay compared to observations (i.e., the simulated monsoon period is too short), while positive values

351    indicate the opposite.

352        For CMIP5, we find systematic errors in the phase of the annual cycle of rainfall. The models are delayed in

353    the onset of summer rainfall over India, the Gulf of Guinea, and the South American Monsoon, with early onset

354    prevalent for the Sahel and the North American Monsoon. The lack of consistency in the phase error across all domains

355    suggests that a ''global'' approach to the study of monsoons may not be sufficient to rectify the regional differences.

356    Rather, regional process studies are necessary for diagnosing the underlying causes of the regionally specific

357    systematic model biases over the different monsoon domains. Assessment of the monsoon fidelity in CMIP6 models

358    using the PMP is in progress.

359

360    *3.6 Cloud feedback and mean-state*

361    Uncertainties in cloud feedback are the primary driver of model-to-model differences in climate sensitivity – the global

362    temperature response to a doubling of atmospheric $CO_2$. Recently, an expert synthesis of several lines of evidence

363    spanning theory, high-resolution models, and observations was conducted to establish quantitative benchmark values

364    (and uncertainty ranges) for several key cloud feedback mechanisms. The assessed feedbacks are those due to changes

365    in high-cloud altitude, tropical marine low-cloud amount, tropical anvil cloud area, land cloud amount, middle latitude

366    marine low-cloud amount, and high latitude low-cloud optical depth. The sum of these six components yields the total

367    assessed cloud feedback, which is part of the overall radiative feedback that fed into the Bayesian calculation of

368    climate sensitivity in Sherwood et al. (2020). Zelinka et al. (2022) estimated these same feedback components in

369    climate models and evaluated them against the expert-judgment values determined in Sherwood et al. (2020),

370    ultimately deriving a root mean square error metric that quantifies the overall match between each model's cloud

371    feedback and those determined through expert judgment.

372        Figure 8 shows the model-simulated values for each individual feedback computed in *amip-p4K* simulations

373    as part of CMIP5 and CMIP6 alongside the expert judgment values. Each model is color-coded by its equilibrium

374    climate sensitivity (determined using *abrupt-4CO2* simulations as described in Zelinka et al., 2020), and the values

375    from an illustrative model (GFDL-CM4) are highlighted. Among the key results apparent from this figure is that

376    models typically underestimate the strength of both positive tropical marine low-cloud feedback and the negative anvil

377    cloud feedback relative to the central expert assessed value. The sum of all six assessed feedback components is

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

380 In addition to evaluating the ability of models to match the assessed cloud feedback components, Zelinka et
381 al. (2022) investigated whether models with less erroneous mean-state clouds tend to have smaller errors in their
382 overall cloud feedback RMSE. This involved computing the mean-state cloud property error metric developed by
383 Klein et al. (2013). This error metric quantifies the spatiotemporal error in climatological cloud properties for clouds
384 with optical depths greater than 3.6, weighted by their net TOA radiative impact. The observational baseline against
385 which the models are compared comes from the ISCCP HGG dataset (Young et al., 2018). Zelinka et al. (2022)
386 showed that models with smaller mean-state cloud errors tend to have stronger but not necessarily better (less
387 erroneous) cloud feedback, which suggests that improving mean-state cloud properties does not guarantee
388 improvement in the cloud response to warming. However, the models with the smallest errors in cloud feedback tend
389 to also have less erroneous mean-state cloud properties, and no models with poor mean-state cloud properties have
390 feedback in good agreement with expert judgment.

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

### 3.7 Precipitation

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

408 The *baseline* metrics gauge the consistency between ESMs and observations, focusing on the holistic set of
409 observed rainfall characteristics (Fig. 9). For example, the spatial distribution of mean state precipitation and seasonal
410 cycle are outcomes of the PMP's Climatology metrics (described in Section 3.1), which provides collective evaluation
411 statistics such as RMSE, standard deviation, and pattern correlation over various domains (e.g., global, NH and SH
412 extratropics, and Tropics, with each domain as a whole, and over land and ocean, in separate). Evaluation of
413 precipitation variability across many timescales with PMP is documented in Ahn et al. (2022); we summarize some
414 of the findings here. The precipitation variability metric measures forced (diurnal and annual cycles) and internal

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

### 3.8 Relating metrics to underlying diagnostics

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

451  PMP's interactive plots have been developed using Bokeh (https://bokeh.org/), a Python data visualization library that

452  enables the creation of interactive plots and applications for web browsers.

453

454  **4 Model Benchmarking**

455  While the PMP originally focused on evaluating multiple models (e.g., Gleckler et al., 2008), in parallel there has

456  been increasing interest from model developers and modeling centers to leverage the PMP to track performance

457  evolution in the model development cycle, as discussed in Gleckler et al. (2016). For example, metrics from the PMP

458  have been used to document performance of ESMs developed in the U.S. DOE Exascale Earth System Model (E3SM;

459  Caldwell et al., 2019; Golaz et al., 2019; Rasch et al., 2019; Hannah et al., 2021; Tang et al., 2021), NOAA

460  Geophysical Fluid Dynamics Laboratory (GFDL; Zhao et al., 2018), Institut Pierre-Simon Laplace (IPSL; Boucher et

461  al., 2020; Planton et al., 2021), National Institute of Meteorological Sciences-Korea Meteorological Administration

462  (NIMS-KMA; Sung et al., 2021), University of California, Los Angeles (Lee et al., 2019b), and the Community

463  Integrated Earth System Model (CIESM) project (Lin et al., 2020).

464  To make the PMP more accessible and useful for modeling groups, efforts are underway to broaden workflow

465  options. Currently, a typical application involves computing a particular class of performance metrics (e.g., mean

466  climate) for all CMIP simulations available via ESGF. To facilitate the ability of modeling groups to routinely use the

467  PMP during their development process, we are working to provide a customized workflow option to run all the PMP

468  metrics more seamlessly on a single model, and to compare these results with a database of PMP results obtained from

469  CMIP simulations (see Code and Data Availability section). Via the PMP-documented and pre-calculated metrics

470  from simulations in the CMIP archive, it is possible to readily incorporate CMIP results into the assessment of new

471  simulations, without retrieving all CMIP simulations and recomputing the results. The resulting quick-look feedback

472  can highlight model improvement (or deterioration) and can assist in determining development priorities or in the

473  selection of a new model version.

474  As an example, here, we show PMP results obtained from GFDL-CM3 from CMIP5 and GFDL-CM4 from

475  CMIP6, for a demonstration of using the Taylor Diagram to compare versions of a given model (Fig. 11). One

476  advantage of the Taylor Diagram is that it collectively represents three statistics (i.e., centered RMSE, standard

477  deviation, and correlation) in a single plot (Taylor, 2001), which synthesizes the performance intercomparison of

478  multiple models (or different versions of a model). In this example, four variables were selected to summarize

479  performance evolution (shown by arrows) in multiple seasons. Except for boreal winter, both model versions are

480  nearly identical in terms of net TOA radiation, however in all seasons the longwave cloud radiative effect is clearly

481  improved in the newer model version. The TOA flux improvements likely contributed to the precipitation

482  improvements, by improving the balances of radiative cooling and latent heating. The improvement in the newer

483  model version is consistent with that documented by Held et al., (2019) and evident via the arrow directions pointing

484  to the observational reference point.

485  Parallel Coordinate Plots can also be used to summarize the comparison of two simulations for their

486  performance. In this section, as an example we demonstrate the comparison of selected metrics: the mean climate,

487  ENSO, and ETMoV (Fig. 12). To facilitate comparison of a subset of models, a few models can be selected and

488 highlighted as connected lines across individual vertical axes on the plot. With the PMP, a common application is to
489 select two versions of the same model to contrast their performance (solid lines) against the backdrop of results from
490 other models, shown as violin plots for the distribution of statistics from other models on each vertical axis. The
491 spatiotemporal RMSE (i.e., temporally averaged spatial RMSE of annual cycle climatology patterns) is used for mean
492 climate as discussed in Section 3.1. The PMP's ENSO metrics that were discussed in Section 3.2 and the RMSE
493 representing total error of ETMoV that were discussed in Section 3.3 are respectively used for ENSO and ETMoV.
494 The plot is simplified from Figure 2 to more efficiently highlight the difference in performance of two GFDL models:
495 GFDL-CM3 and GFDL-CM4. Each vertical axis indicates performance for each metric defined for climatology of
496 variables (Fig. 12a), ENSO characteristics (Fig. 12b), or interannual variability mode obtained from seasonal or
497 monthly averaged time series (Fig. 12c). In this example, it is shown that GFDL-CM4 is superior to GFDL-CM3 for
498 most cases across selected metrics (downward arrows in green) while inferior for a few cases (upward arrows in red)
499 — consistent with previous findings (Held et al., 2019; Planton et al., 2021; Chen et al., 2021). Such applications of
500 the Parallel Coordinate Plot can enable quick overall assessment and tracking of the ESM performance evolution
501 during its development cycle. More examples showing other models are available in the Supplementary material (Figs.
502 S1 to S3).

503 Note that there have been efforts to coalesce objective model evaluation concepts used in the research
504 community (e.g., Knutti et al., 2010), however as the field continues to evolve rapidly, definitions are still being
505 finessed, and there is room for the community to further advance well-established metrics. Via the PMP, we produce
506 hundreds of summary statistics, but it will not be surprising if only a subset of them might be considered as viable
507 candidate metrics for more practical routine performance evaluations.

508

509 **5 Discussion**

510 Given the critical role ESMs play in our efforts to understand a changing climate, scientists involved in the analysis
511 of ESM simulations have been compelled to improve the process of model evaluation. Current progress towards
512 systematic model evaluation remains dynamic, with evolving approaches and many independent paths being pursued.
513 This has resulted in the development of diversified model evaluation software packages. For example, ESMValTool
514 (Eyring et al., 2016, 2019, 2020; Righi et al., 2020) is a comprehensive package led by a European core development
515 team that has been used for numerous applications including producing model evaluation plots in Chapter 3 of the
516 IPCC's AR6 Working Group 1 Assessment (Eyring et al., 2021). The Model Diagnostics Task Force (MDTF)
517 Diagnostics package, led by NOAA, focuses on process-oriented diagnostics (Maloney et al., 2019; Neelin et al.,
518 2023). The International Land Model Benchmarking (ILAMB) Software System (Collier et al., 2018) led by Oak
519 Ridge National Laboratory provides land surface and carbon cycle metrics with key state-ot-the art observational
520 products, and similarly, the International Ocean Model Benchmarking (IOMB) Software System (Fu et al., 2022)
521 focuses on surface and upper ocean biogeochemical variables. The Climate Variability Diagnostics Package (CVDP;
522 Phillips et al., 2014; Fasullo et al., 2020) developed at NCAR provides diagnosis of climate modes of variability.
523 Analyzing Scales of Precipitation (ASoP; Klingaman et al., 2017; Martin et al., 2017; Ordonez et al., 2021) focuses
524 on analyzing precipitation scales across space and time. In parallel, the regional climate community also has actively

525    developed metrics packages such as the Regional Climate Model Evaluation System (RCMES; Lee et al., 2018a;
526    Whitehall et al. 2012). Separately, a few climate modeling centers have developed their own model evaluation
527    packages to assist in their in-house ESM development, e.g., the E3SM Diags (Zhang et al., 2022). There also have
528    been other efforts to enhance the usability of in-situ and field campaign observations in ESM evaluations, such as
529    Atmospheric Radiation Measurement (ARM) GCM Diag (Zhang et al., 2018, 2020) and Earth System Model Aerosol–
530    Cloud Diagnostics (ESMAC Diags; Tang et al., 2022, 2023).

531         The model evaluation packages currently being advanced within the ESM research community all have their
532    own technical approaches and scientific priorities. We believe that this diversity has made, and will continue to make,
533    the model evaluation process even more comprehensive and successful. The fact that there is some overlap in a few
534    cases is advantageous because it enables the cross-verification of results, which is particularly useful in the more
535    complex analyses. Despite the advantages, having no single best or widely accepted approach for the community to
536    follow, does introduce complexity to the coordination of model evaluation. To facilitate collective usages of individual
537    evaluation tools, the CMEC has initiated the development of a unified code base that technically coordinates the
538    operation of distinct but complementary tools (Ordonez et al. 2021). Currently, the PMP, ILAMB, MDTF and ASoP
539    have become CMEC-compliant by adopting the common interface standards that define how evaluation tools interact
540    with observational data and climate model output. We expect that CMEC can also help the model evaluation
541    community to establish standards for archiving the metrics output, much as the community did for the conventions to
542    describe climate model data (e.g., CMIP application of CF Metadata Conventions [http://cfconventions.org/]; Hassell
543    et al., 2017; Eaton et al., 2022).

544         It is worth noting that the comprehensive database of PMP results offers a resource for exploring the range
545    of structural errors in CMIP class models and their interrelationships. For example, examination of cross-metric
546    relationships between mean-state and variability biases can shed additional light on the propagation of errors (e.g.,
547    Kang et al., 2020; Lee et al., 2021b). There continues to be interest in ranking models for specific applications (e.g.,
548    Ashfaq et al., 2022; Goldenson et al., 2023; Longmate et al., 2023; Papalexiou et al., 2020) or to "move beyond one
549    model one vote" in multi-model analysis to reduce uncertainties in the spread of multi-model projections (e.g., Knutti,
550    2010; Knutti et al., 2017; Sanderson et al., 2017; Herger et al., 2018; Hausfather et al., 2022; Merrifield et al., 2023).
551    While we acknowledge potential interests in using the results of the PMP or equivalent to rank models or identify
552    performance outliers (e.g., Sanderson and Wehner, 2017), we believe the many challenges associated with model
553    weighting are application dependent, and thus leave it up to users of the PMP to make those judgments.

554

**6 Summary and Future Directions**

556    The PMP has provided quasi-operational ESM evaluation capabilities that can be rapidly deployed to objectively
557    summarize a diverse suite of model behavior with results made publically available. This can be of value in the
558    assessment of community intercomparisons like CMIP, the evaluation of large ensembles, or the model development
559    process. By documenting objective performance summaries produced by the PMP and making them available via
560    detailed version control, additional research is made possible beyond the baseline model evaluation, model
561    intercomparison, and benchmarking. The outcomes of PMP's calculations applied to the CMIP archive culminate in

562    the PCMDI Simulation Summary (https://pcmdi.llnl.gov/metrics/). This summary serves as a comprehensive

563    repository of PMP outputs, visually capturing the outcomes of objective model-to-observation comparisons. Special

564    attention is dedicated to the most recent ensemble of models contributing to CMIP6. By offering a comprehensive

565    assessment of simulated climate, its variability modes, and characteristics of precipitation in ESMs, the PMP

566    framework equips model developers with quantifiable benchmarks to validate and enhance model performance.

567          With the growing interest in augmenting the suite of metrics within PMP that reflects an evolving landscape

568    of evaluation needs, continual efforts are channeled into expanding the scope of the PMP. For example, in coordination

569    with the World Meteorological Organization (WMO)'s WGNE MJO Task Force, additional candidate MJO metrics

570    for PMP inclusion have been identified to facilitate more comprehensive assessments of the MJO. Implementation of

571    metrics for MJO amplitude, periodicity, and structure into the PMP is planned. The ongoing collaboration with NCAR

572    aims to incorporate metrics related to the upper atmosphere, specifically the Quasi-Biennial Oscillation (QBO) and

573    QBO-MJO metrics (e.g. Kim et al., 2020). We also have plans to grow the scope of PMP beyond its traditional

574    atmospheric realm to include domains like the ocean and Arctic regions through collaboration with the U.S. DOE's

575    project entitled High Latitude Application and Testing of ESMs (HiLAT, https://www.hilat.org/). This dimension of

576    evaluation holds promise in offering deeper insights into model performance.

577          In addition to the scientific challenges associated with diversifying objective summaries of model

578    performance, there are numerous potential areas to advance accompanying technologies, in large part related to the

579    rapidly evolving set of open-source tools and methods available to scientists. We expect that the current ongoing PMP

580    code modernization effort to fully adapt the xCDAT will potentially galvanize greater community involvement. We

581    will continue to maintain robust rigorousness in the calculation of statistics for the PMP metrics by staying tuned with

582    the latest progress in the field, such as implementing the method for more rigorous conservation in horizontal

583    interpolation (Taylor, 2023). To improve clarity of key deliverable messages from multivariate data of PMP's metrics

584    obtained from comprehensive ESM archives, we will consider implementing the advances in the high-dimensional

585    data visualization field, such as the circular plot discussed in Lee et al. (2018b) and variations of Parallel Coordinate

586    Plots proposed by Hassan et al. (2019) and Lu et al. (2020).

587          Looking ahead, the PMP framework is also poised to contribute to high-resolution climate modeling

588    communities, notably the High Resolution Model Intercomparison Project (HighResMIP; Haarsma et al., 2016) and

589    the DYnamics of the Atmospheric general circulation Modeled On Non-hydrostatic Domains (DYAMOND; Stevens

590    et al., 2019). This motivates developments of specialized metrics for high-resolution models, which demonstrate the

591    features that high-resolution models have enabled. Potential avenue of exploration involves leveraging Machine

592    Learning (ML) techniques, considering the examined applicability of ML and other state-of-the art data science

593    techniques being used for process-oriented ESM evaluation works (e.g., Nowack et al., 2020; Labe and Barnes, 2022;

594    Dalelane et al., 2023). Applications of ML detections, such as for storms using TempestExtremes (Ullrich and

595    Zarzycki 2017; Ullrich et al., 2021) and fronts (e.g, Biard and Kunkel, 2019), can enable additional specialized storm

596    metrics for high resolution simulations. For convection permitting models, yet more storm metrics can be applied such

597    as Mesoscale convective systems. Into the PMP, we currently have plans to implement atmospheric blocking metrics

598    that were developed through the collaboration of Colorado State University and the PCMDI (Valkonen et al., in prep),

and Atmospheric River detection metrics that are currently under development at LLNL. Both of these metrics suites were developed using the pattern detection capabilities in the latest TempestExtremes (Ullrich et al., 2021). This application of the PMP aligns with a broader plan for regional expansion, with a deliberate emphasis on processes intrinsic to specific regions.

We anticipate that the PMP will continue to play a crucial role in benchmarking ESMs in the future. Improvements in PMP, coupled with advancements in projects within the MIP community, will significantly contribute to assessing the evolving performance of ESMs including via the collaboration with the CMIP Benchmarking Task Team. Enhancements in version control and transparency within obs4MIPs are poised to enhance the provenance and reproducibility of PMP results, thereby strengthening the foundation for rigorous and repeatable performance benchmarking. The PMP's collaboration with the CMIP Forcing Task Team, through the Input4MIPs (Durack et al., 2018) and the CMIP6Plus projects, will further expand the utility of performance metrics in identifying problems associated with the forcing dataset and their application and use in reproducing the observed record of historical climate. Furthermore, as ESMs advance towards more operationalized configurations to meet the demands of decision-making processes (Jakob et al., 2023), the PMP holds significant potential to provide interoperable ESM evaluation and benchmarking capabilities to the community.

**Author Contributions**

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

**Code and Data Availability**

The source code of PMP (Lee et al., 2023b) is available as an open-source Python package: https://github.com/PCMDI/pcmdi_metrics (last access: 21 November 2023) with versions archived on Zenodo DOI: https://doi.org/10.5281/zenodo.592790 (last access: 21 November 2023). The PMP results database (Lee et al., 2023a) that includes calculated metrics is available on the GitHub repository at https://github.com/PCMDI/pcmdi_metrics_results_archive (last access: 21 November 2023) with versions archived on Zenodo DOI: https://doi.org/10.5281/zenodo.10181201. The interactive visualizations of the PMP results are available on the PCMDI website at https://pcmdi.llnl.gov/metrics (last access: 21 November 2023). The CMIP5 and CMIP6 model outputs and obs4MIPs datasets used in this paper are available via the Earth System Grid Federation at https://esgf-node.llnl.gov/ (last access: 21 November 2023).

**Competing interests**

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

**References**

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

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

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

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

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

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

679   Bayr, T., Wengel, C., Latif, M., Dommenget, D., Lübbecke, J., and Park, W.: Error compensation of ENSO
680         atmospheric feedbacks in climate models and its influence on simulated ENSO dynamics, Climate Dynamics,
681         53, 155–172, https://doi.org/10.1007/s00382-018-4575-7, 2019.

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

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

686   Boucher, O., Servonnat, J., Albright, A. L., Aumont, O., Balkanski, Y., Bastrikov, V., Bekki, S., Bonnet, R., Bony,
687         S., Bopp, L., Braconnot, P., Brockmann, P., Cadule, P., Caubel, A., Cheruy, F., Codron, F., Cozic, A., Cugnet,
688         D., D'Andrea, F., Davini, P., De Lavergne, C., Denvil, S., Deshayes, J., Devilliers, M., Ducharne, A.,
689         Dufresne, J. L., Dupont, E., Ethé, C., Fairhead, L., Falletti, L., Flavoni, S., Foujols, M. A., Gardoll, S.,
690         Gastineau, G., Ghattas, J., Grandpeix, J. Y., Guenet, B., Lionel, E. G., Guilyardi, E., Guimberteau, M.,
691         Hauglustaine, D., Hourdin, F., Idelkadi, A., Joussaume, S., Kageyama, M., Khodri, M., Krinner, G., Lebas,
692         N., Levavasseur, G., Lévy, C., Li, L., Lott, F., Lurton, T., Luyssaert, S., Madec, G., Madeleine, J.-B.,
693         Maignan, F., Marchand, M., Marti, O., Mellul, L., Meurdesoif, Y., Mignot, J., Musat, I., Ottlé, C., Peylin, P.,
694         Planton, Y., Polcher, J., Rio, C., Rochetin, N., Rousset, C., Sepulchre, P., Sima, A., Swingedouw, D.,
695         Thiéblemont, R., Traore, A. K., Vancoppenolle, M., Vial, J., Vialard, J., Viovy, N., and Vuichard, N.:
696         Presentation and evaluation of the IPSL-CM6A-LR Climate Model, Journal of Advances in Modeling Earth
697         Systems, 12, https://doi.org/10.1029/2019ms002010, 2020.

698   Caldwell, P., Mametjanov, A., Tang, Q., Van Roekel, L., Golaz, J.-C., Lin, W., Bader, D. C., Keen, N. D., Feng, Y.,
699         Jacob, R., Maltrud, M., Roberts, A., Taylor, M. A., Veneziani, M., Wang, H., Wolfe, J. D., Balaguru, K.,
700         Cameron-Smith, P. J., Dong, L., Klein, S. A., Leung, L. R., Li, H., Li, Q., Liu, X., Neale, R., Pinheiro, M.
701         C., Qian, Y., Ullrich, P. A., Xie, S., Yang, Y., Zhang, Y., Zhang, K., and Zhou, T.: The DOE E3SM Coupled
702         Model Version 1: description and results at high resolution, Journal of Advances in Modeling Earth Systems,
703         11, 4095–4146, https://doi.org/10.1029/2019ms001870, 2019.

704   Cannon, A. J.: Reductions in daily continental-scale atmospheric circulation biases between generations of global
705         climate   models:   CMIP5   to   CMIP6,   Environmental   Research   Letters,   15,   064006,
706         https://doi.org/10.1088/1748-9326/ab7e4f, 2020.

707    Chen, H.-C., Jin, F.-F., Zhao, S., Wittenberg, A. T., and Xie, S.: ENSO dynamics in the E3SM-1-0, CESM2, and
708         GFDL-CM4 climate models, Journal of Climate, 34, 9365-9384, https://doi.org/10.1175/JCLI-D-21-0355.1,
709         2021.

710    Choi, J. H. and Son, S.-W.: Seasonal-to-decadal prediction of El Niño–Southern Oscillation and Pacific Decadal
711         Oscillation, Npj Climate and Atmospheric Science, 5, https://doi.org/10.1038/s41612-022-00251-9, 2022.

712    Collier, N., Hoffman, F. M., Lawrence, D. M., Keppel-Aleks, G., Koven, C. D., Riley, W. J., Mu, M., and Randerson,
713         J. T.: The International Land Model Benchmarking (ILAMB) System: Design, theory, and implementation,
714         Journal of Advances in Modeling Earth Systems, 10, 2731–2754, https://doi.org/10.1029/2018ms001354,
715         2018.

716    Covey, C., AchutaRao, K., Cubasch, U., Jones, P., Lambert, S. J., Mann, M., Phillips, T. J., and Taylor, K. E.: An
717         overview of results from the Coupled Model Intercomparison Project, Global and Planetary Change, 37, 103–
718         133, https://doi.org/10.1016/s0921-8181(02)00193-5, 2003.

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

722    Crockford, D.: The application/json media type for javascript object notation (json) (No. rfc4627), https://www.rfc-
723         editor.org/rfc/pdfrfc/rfc4627.txt.pdf (last access: 6 November 2023), 2006.

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

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

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

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

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

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

741    Eaton, B., Gregory, J., Drach, B., Taylor, K., Hankin, S., Blower, J., Caron, J., Signell, R., Bentley, P., Rappa, G.,
742         Höck, H., Pamment, A., Juckes, M., Raspaud, M., Horne, R., Whiteaker, T., Blodgett, D., Zender, C., Lee,
743         D., Hassell, D., Snow, A. D., Kölling, T., Allured, D., Jelenak, A., Soerensen, A. M., Gaultier, L., Herlédan,

744    S.: NetCDF Climate and Forecast (CF) Meta-data Conventions V1.10, available at:
745    http://cfconventions.org/Data/cf-conventions/cf-conventions-1.10/cf-conventions.html (last access: 6
746    November 2023), 2022.

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

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

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

763  Eyring, V., Bock, L., Lauer, A., Righi, M., Schlund, M., Andela, B., Arnone, E., Bellprat, O., Brötz, B., Caron, L.-P.,
764    Carvalhais, N., Cionni, I., Cortesi, N., Crezee, B., Davin, E. L., Davini, P., Debeire, K., De Mora, L., Deser,
765    C., Docquier, D., Earnshaw, P., Ehbrecht, C., Gier, B. K., Gonzalez-Reviriego, N., Goodman, P. J.,
766    Hagemann, S., Hardiman, S. C., Hassler, B., Hunter, A., Kadow, C., Kindermann, S., Koirala, S., Koldunov,
767    N., Lejeune, Q., Lembo, V., Lovato, T., Lucarini, V., Massonnet, F., Müller, B., Pandde, A., Pérez-Zanón,
768    N., Phillips, A. S., Predoi, V., Russell, J. L., Sellar, A., Serva, F., Stacke, T., Swaminathan, R., Torralba, V.,
769    Vegas-Regidor, J., Von Hardenberg, J., Weigel, K., and Zimmermann, K.: Earth System Model Evaluation
770    Tool (ESMValTool) v2.0 – an extended set of large-scale diagnostics for quasi-operational and
771    comprehensive evaluation of Earth system models in CMIP, Geoscientific Model Development, 13, 3383–
772    3438, https://doi.org/10.5194/gmd-13-3383-2020, 2020.

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

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

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

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

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

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

794 Gates, W.L.: AN AMS continuing series: Global CHANGE–AMIP: The atmospheric model intercomparison project,
795     Bulletin of the American Meteorological Society, 73, 1962-1970, 1992.

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

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

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

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

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

808 Golaz, J.-C., Caldwell, P., Van Roekel, L., Petersen, M. R., Tang, Q., Wolfe, J. D., Abeshu, G. W., Anantharaj, V.,
809     Asay-Davis, X., Bader, D. C., Baldwin, S., Bisht, G., Bogenschutz, P., Branstetter, M. L., Brunke, M. A.,
810     Brus, S., Burrows, S. M., Cameron-Smith, P. J., Donahue, A. S., Deakin, M., Easter, R. C., Evans, K. J.,
811     Feng, Y., Flanner, M., Foucar, J. G., Fyke, J., Griffin, B. M., Hannay, C., Harrop, B. E., Hoffman, M. J.,
812     Hunke, E., Jacob, R., Jacobsen, D. W., Jeffery, N., Jones, P. W., Keen, N. D., Klein, S. A., Larson, V. E.,
813     Leung, L. R., Li, H. Y., Lin, W., Lipscomb, W. H., Lun, P., Mahajan, S., Maltrud, M., Mametjanov, A.,
814     McClean, J. L., McCoy, R., Neale, R., Price, S., Qian, Y., Rasch, P. J., Eyre, J. E. J. R., Riley, W. J., Ringler,
815     T. D., Roberts, A., Roesler, E. L., Salinger, A. G., Shaheen, Z., Shi, X., Singh, B., Tang, J., Taylor, M. A.,
816     Thornton, P. E., Turner, A. K., Veneziani, M., Wan, H., Wang, H., Wang, S., Williams, D. N., Wolfram, P.
817     J., Worley, P. H., Xie, S., Yang, Y., Yoon, J., Zelinka, M. D., Zender, C. S., Zeng, X., Zhang, C., Zhang, K.,

818        Zhang, Y., Zheng, X., Zhou, T., and Zhu, Q.: The DOE E3SM Coupled Model Version 1: Overview and
819        evaluation at standard resolution, Journal of Advances in Modeling Earth Systems, 11, 2089–2129,
820        https://doi.org/10.1029/2018ms001603, 2019.

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

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

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

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

837  Hintze, J. L., and Nelson, R. D.: Violin plots: A box plot-density trace synergism, The American Statistician, 52, 181–
838        184, https://doi.org/10.1080/00031305.1998.10480559, 1998.

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

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

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

848  Hausfather, Z., Marvel, K., Schmidt, G. A., Nielsen-Gammon, J. W. and Zelinka, M.: Climate simulations: Recognize
849        the 'hot model' problem, Nature, 605, 26-29, https://doi.org/10.1038/d41586-022-01192-2, 2022.

850  Held, I. M., Guo, H., Adcroft, A., Dunne, J. P., Horowitz, L. W., Krasting, J., Shevliakova, E., Winton, M., Zhao, M.,
851        Bushuk, M., Wittenberg, A. T., and coauthors: Structure and performance of GFDL's CM4. 0 climate model,
852        Journal of Advances in Modeling Earth Systems, 11, 3691-3727, https://doi.org/10.1029/2019MS001829,
853        2019.

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

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

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

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

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

868    Inselberg, A.: Multidimensional detective, in: Proceedings of IEEE Symposium on Information Visualization, 100–
869          107, https://doi.org/10.1109/INFVIS.1997.636793, 1997.

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

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

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

878    Jakob, C., Gettelman, A. and Pitman, A.: The need to operationalize climate modelling, Nat. Clim. Chang. 13, 1158–
879          1160, https://doi.org/10.1038/s41558-023-01849-4, 2023.

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

883    Jun, S.-Y., Kim, J.-H., Choi, J. H., Kim, S.-J., Kim, B.-M., and An, S.-I.: The internal origin of the west-east
884          asymmetry of Antarctic climate change, Science Advances, 6, https://doi.org/10.1126/sciadv.aaz1490, 2020.

885    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
886          in CESM2 ensemble simulation, Geophysical Research Letters, 47, https://doi.org/10.1029/2020gl089824,
887          2020.

888    Kim, D., Sperber, K. R., Stern, W., Waliser, D. E., Kang, I. S., Maloney, E. D., Wang, W., Weickmann, K. M.,
889          Benedict, J. J., Khairoutdinov, M., Lee, M.-I., Neale, R., Suarez, M. J., Thayer-Calder, K., and Zhang, G.:

890    Application of MJO simulation diagnostics to climate models, Journal of Climate, 22, 6413–6436,
891        https://doi.org/10.1175/2009jcli3063.1, 2009.
892  Kim, H., Caron, J. M., Richter, J. H. and Simpson, I. R.: The lack of QBO-MJO connection in CMIP6 models,
893        Geophysical Research Letters, 47, e2020GL087295, https://doi.org/10.1029/2020GL087295, 2020.
894  Kim, Y.-H., Min, S.-K., Zhang, X., Sillmann, J., and Sandstad, M.: Evaluation of the CMIP6 multi-model ensemble
895        for    climate    extreme    indices,    Weather    and    Climate    Extremes,    29,    100269,
896        https://doi.org/10.1016/j.wace.2020.100269, 2020.
897  Klein, S. A., Zhang, Y., Zelinka, M. D., Pincus, R., Boyle, J., and Gleckler, P. J.: Are climate model simulations of
898        clouds improving? An evaluation using the ISCCP simulator, Journal of Geophysical Research:
899        Atmospheres, 118, 1329–1342, https://doi.org/10.1002/jgrd.50141, 2013.
900  Klingaman, N. P., Martin, G., and Moise, A.: ASoP (v1.0): a set of methods for analyzing scales of precipitation in
901        general circulation models, Geoscientific Model Development, 10, 57–83, https://doi.org/10.5194/gmd-10-
902        57-2017, 2017.
903  Knutti, R.: The end of model democracy? Climatic Change, 102, 395–404, https://doi.org/10.1007/s10584-010-9800-
904        2, 2010.
905  Knutti, R., Sedláček, J., Sanderson, B. M., Lorenz, R., Fischer, E. M., and Eyring, V.: A climate model projection
906        weighting scheme accounting for performance and interdependence, Geophysical Research Letters,
907        https://doi.org/10.1002/2016gl072012, 2017.
908  Knutti, R., Abramowitz, G., Collins, M., Eyring, V., Gleckler, P. J., Hewitson, B., and Mearns, L.: Good Practice
909        Guidance Paper on Assessing and Combining Multi Model Climate Projections, in: Meeting Report of the
910        Intergovernmental Panel on Climate Change Expert Meeting on Assessing and Combining Multi Model
911        Climate Projections, edited by Stocker, T. F., Qin, D., Plattner, G.-K., Tignor, M., and Midgley, P. M., IPCC
912        Working Group I Technical Support Unit, University of Bern, Bern, Switzerland, 2010.
913  Labe, Z. M. and Barnes, E. A.: Comparison of Climate Model Large Ensembles With Observations in the Arctic Using
914        Simple    Neural    Networks,    Earth    and    Space    Science,    e2022EA002348,
915        https://doi.org/10.1029/2022EA002348, 2022.
916  Lambert, S. J. and Boer, G. J.: CMIP1 evaluation and intercomparison of coupled climate models, Climate Dynamics,
917        17, 83-106, https://doi.org/10.1007/PL00013736, 2001.
918  Lee, H., Goodman, A., McGibbney, L. J., Waliser, D. E., Kim, J., Loikith, P. C., Gibson, P. B., and Massoud, E.:
919        Regional Climate Model Evaluation System powered by Apache Open Climate Workbench v1.3.0: an
920        enabling tool for facilitating regional climate studies, Geoscientific Model Development, 11, 4435–4449,
921        https://doi.org/10.5194/gmd-11-4435-2018, 2018a.
922  Lee, J., Ahn, M.-S., Ordonez, A., Gleckler, P., and Ullrich, P.: PCMDI/pcmdi_metrics_results_archive, Zenodo [data],
923        https://doi.org/10.5281/zenodo.10181201, 2023a.
924  Lee, J., Gleckler, P., Ordonez, A., Ahn, M.-S., Ullrich, P., Tom, V., Jason, B., Charles, D., Durack, P., Shaheen, Z.,
925        Muryanto, L., Painter, J., and Krasting, J.: PCMDI/pcmdi_metrics: PMP Version 3.1.1, Zenodo [code],
926        https://doi.org/10.5281/zenodo.592790, 2023b.

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

980 McPhaden, M. J., Zebiak, S. E., and Glantz, M. H.: ENSO as an integrating concept in Earth Science, Science, 314,
981      1740–1745, https://doi.org/10.1126/science.1132588, 2006.

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

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

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

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

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

994 Merrifield, A., Brunner, L., Lorenz, R., Humphrey, V., and Knutti, R.: Climate model Selection by Independence,
995      Performance, and Spread (ClimSIPS v1.0.1) for regional applications, Geoscientific Model Development,
996      16, 4715–4747, https://doi.org/10.5194/gmd-16-4715-2023, 2023.

997 Neelin, J. D., Krasting, J. P., Radhakrishnan, A., Liptak, J., Jackson, T. J., Ming, Y., Dong, W., Gettelman, A.,
998      Coleman, D., Maloney, E. D., Wing, A. A., Kuo, Y. H., Ahmed, F., Ullrich, P. A., Bitz, C. M., Neale, R.,
999      Ordonez, A., and Maroon, E.: Process-oriented diagnostics: principles, practice, community development

1000    and common standards, Bulletin of the American Meteorological Society, https://doi.org/10.1175/bams-d-
1001        21-0268.1, 2023.

1002    Nowack, P., Runge, J., Eyring, V., and Haigh, J. D.: Causal networks for climate model evaluation and constrained
1003        projections, Nature Communications, 11, https://doi.org/10.1038/s41467-020-15195-y, 2020.

1004    Orbe, C., Van Roekel, L., Adames, Á. F., Dezfuli, A., Fasullo, J. T., Gleckler, P. J., Lee, J., Li, W., Nazarenko, L.,
1005        Schmidt, G. A., Sperber, K. R., and Zhao, M.: Representation of modes of variability in six U.S. climate
1006        models, Journal of Climate, 33, 7591–7617, https://doi.org/10.1175/jcli-d-19-0956.1, 2020.

1007    Ordonez, A. C., Klingaman, N. P., and Martin, G.: Analysing scales of precipitation, OSTI OAI (U.S. Department of
1008        Energy Office of Scientific and Technical Information), https://doi.org/10.11578/dc.20211029.5, 2021.

1009    Papalexiou, S. M., Rajulapati, C. R., Clark, M. P., and Lehner, F.: Robustness of CMIP6 historical global mean
1010        temperature simulations: Trends, long-term persistence, autocorrelation, and distributional shape, Earth's
1011        Future, 8, e2020EF001667, https://doi.org/10.1029/2020EF001667, 2020.

1012    Pascoe, C., Lawrence, B. N., Guilyardi, E., Juckes, M., and Taylor, K. E.: Documenting numerical experiments in
1013        support of the Coupled Model Intercomparison Project Phase 6 (CMIP6), Geosci. Model Dev., 13, 2149–
1014        2167, https://doi.org/10.5194/gmd-13-2149-2020, 2020.

1015    Pendergrass, A. G., Gleckler, P. J., Leung, L. R., and Jakob, C.: Benchmarking simulated precipitation in earth system
1016        models, Bulletin of the American Meteorological Society, 101, E814–E816, https://doi.org/10.1175/bams-d-
1017        19-0318.1, 2020.

1018    Phillips, A. S., Deser, C., and Fasullo, J. T.: Evaluating modes of variability in climate models, Eos, Transactions
1019        American Geophysical Union, 95, 453–455, https://doi.org/10.1002/2014eo490002, 2014.

1020    Planton, Y., Guilyardi, E., Wittenberg, A. T., Lee, J., Gleckler, P. J., Bayr, T., McGregor, S., McPhaden, M. J., Power,
1021        S. B., Roehrig, R., Vialard, J., and Voldoire, A.: Evaluating Climate Models with the CLIVAR 2020 ENSO
1022        Metrics Package, Bulletin of the American Meteorological Society, 102, E193–E217,
1023        https://doi.org/10.1175/bams-d-19-0337.1, 2021.

1024    Potter, G. L., Bader, D. C., Riches, M., Bamzai, A. and Joseph, R.: Celebrating two decades of the Program for Climate
1025        Model Diagnosis and Intercomparison. Bulletin of the American Meteorological Society, 92, 629-631,
1026        https://doi.org/10.1175/2011BAMS3018.1, 2011.

1027    Qin, Y., Zelinka, M. D., and Klein, S. A.: On the Correspondence Between Atmosphere-Only and Coupled
1028        Simulations for Radiative Feedbacks and Forcing From $CO_2$, Journal of Geophysical Research:
1029        Atmospheres, 127, https://doi.org/10.1029/2021jd035460, 2022.

1030    Randall, D.A., Wood, R.A., Bony, S., Colman, R., Fichefet, T., Fyfe, J., Kattsov, V., Pitman, A., Shukla, J., Srinivasan,
1031        J. and Stouffer, R.J.: Climate models and their evaluation. In Climate change 2007: The physical science
1032        basis. Contribution of Working Group I to the Fourth Assessment Report of the IPCC (FAR), 589-662,
1033        Cambridge University Press, 2007.

1034    Randerson, J. T., Hoffman, F. M., Thornton, P. E., Mahowald, N. M., Lindsay, K., Lee, Y.-H., Nevison, C. D., Doney,
1035        S. C., Bonan, G. B., Stöckli, R., Covey, C., Running, S. W., and Fung, I.: Systematic assessment of terrestrial

1036    biogeochemistry in coupled climate-carbon models, Global Change Biology, 15, 2462–2484,
1037        https://doi.org/10.1111/j.1365-2486.2009.01912.x, 2009.

1038    Rasch, P. J., Xie, S., Ma, P.-L., Lin, W., Wang, H., Tang, Q., Burrows, S. M., Caldwell, P., Zhang, K., Easter, R. C.,
1039        Cameron-Smith, P. J., Singh, B., Wan, H., Golaz, J.-C., Harrop, B. E., Roesler, E. L., Bacmeister, J. T.,
1040        Larson, V. E., Evans, K. J., Qian, Y., Taylor, M. A., Leung, L. R., Zhang, Y., Brent, L., Branstetter, M. L.,
1041        Hannay, C., Mahajan, S., Mametjanov, A., Neale, R., Richter, J. H., Yoon, J.-H., Zender, C. S., Bader, D. C.,
1042        Flanner, M., Foucar, J. G., Jacob, R., Keen, N. D., Klein, S. A., Liu, X., Salinger, A. G., Shrivastava, M., and
1043        Yang, Y.: An overview of the atmospheric component of the Energy Exascale Earth System model, Journal
1044        of Advances in Modeling Earth Systems, 11, 2377–2411, https://doi.org/10.1029/2019ms001629, 2019.

1045    Reichler, T. and Kim, J.: How well do coupled models simulate today's climate?, Bulletin of the American
1046        Meteorological Society, 89, 303–312, https://doi.org/10.1175/bams-89-3-303, 2008.

1047    Righi, M., Andela, B., Eyring, V., Lauer, A., Predoi, V., Schlund, M., Vegas-Regidor, J., Bock, L., Brötz, B., De
1048        Mora, L., Diblen, F., Dreyer, L., Drost, N., Earnshaw, P., Hassler, B., Koldunov, N., Little, B., Tomas, S. L.,
1049        and Zimmermann, K.: Earth System Model Evaluation Tool (ESMValTool) v2.0 – technical overview,
1050        Geoscientific Model Development, 13, 1179–1199, https://doi.org/10.5194/gmd-13-1179-2020, 2020.

1051    Sanderson, B. M. and Wehner, M. F.: Weighting strategy for the Fourth National Climate Assessment, in: Climate
1052        Science Special Report: Fourth National Climate Assessment, Volume I, edited by Wuebbles, D. J., Fahey,
1053        D. W., Hibbard, K. A., Dokken, D. J., Stewart, B. C., and Maycock, T.K., U.S. Global Change Research
1054        Program, Washington, DC, USA, 436-442, https://doi.org/10.7930/J06T0JS3, 2017.

1055    Sanderson, B. M., Wehner, M., and Knutti, R.: Skill and independence weighting for multi-model assessments,
1056        Geoscientific Model Development, 10, 2379–2395, https://doi.org/10.5194/gmd-10-2379-2017, 2017.

1057    Sherwood, S. C., Webb, M. J., Annan, J. D., Armour, K. C., Forster, P. M., Hargreaves, J. C., Hegerl, G. C., Klein, S.
1058        A., Marvel, K., Rohling, E. J., Watanabe, M., Andrews, T., Braconnot, P., Bretherton, C. S., Foster, G. L.,
1059        Hausfather, Z., Von Der Heydt, A. S., Knutti, R., Mauritsen, T., Norris, J. R., Proistosescu, C., Rugenstein,
1060        M., Schmidt, G. A., Tokarska, K., and Zelinka, M. D.: An assessment of Earth's climate sensitivity using
1061        multiple lines of evidence, Reviews of Geophysics, 58, https://doi.org/10.1029/2019rg000678, 2020.

1062    Sillmann, J., Kharin, V. V., Zhang, X., Zwiers, F. W., and Bronaugh, D.: Climate extremes indices in the CMIP5
1063        multimodel ensemble: Part 1. Model evaluation in the present climate, Journal of Geophysical Research:
1064        Atmospheres, 118, 1716–1733, https://doi.org/10.1002/jgrd.50203, 2013.

1065    Sperber, K. R.: Madden-Julian variability in NCAR CAM2.0 and CCSM2.0, Clim Dyn 23, 259–278,
1066        https://doi.org/10.1007/s00382-004-0447-4, 2004.

1067    Sperber, K. R., Annamalai, H., Kang, I.-S., Kitoh, A., Moise, A., Turner, A., Wang, B., and Zhou, T.: The Asian
1068        summer monsoon: an intercomparison of CMIP5 vs. CMIP3 simulation of the late 20th century. Clim Dyn
1069        41, 2711-2744, https://doi.org/10.1007/s00382-012-1607-6, 2013.

1070    Sperber K. R., Gualdi, S., Legutke, S., Gayler, V.: The Madden–Julian oscillation in ECHAM4 coupled and uncoupled
1071        general circulation models, Clim Dyn 25, 117–140, https://doi.org/10.1007/s00382-005-0026-3, 2005.

1072  Srivastava, A., Grotjahn, R., and Ullrich, P. A.: Evaluation of historical CMIP6 model simulations of extreme
1073      precipitation over contiguous US regions, Weather and Climate Extremes, 29, 100268,
1074      https://doi.org/10.1016/j.wace.2020.100268, 2020.

1075  Steed, C. A., Shipman, G., Thornton, P., Ricciuto, D., Erickson, D. and Branstetter, M.: Practical application of parallel
1076      coordinates for climate model analysis, Procedia Computer Science, 9, 877-886,
1077      https://doi.org/10.1016/j.procs.2012.04.094, 2012.

1078  Stevens, B., Satoh, M., Auger, L., Biercamp, J., Bretherton, C. S., Chen, X., Düben, P., Judt, F., Khairoutdinov, M.,
1079      Klocke, D., Kodama, C., Kornblueh, L., Lin, S.-J., Neumann, P., Putman, W. M., Röber, N., Shibuya, R.,
1080      Vanniere, B., Vidale, P. L., Wedi, N., and Zhou, L.: DYAMOND: the DYnamics of the Atmospheric general
1081      circulation Modeled On Non-hydrostatic Domains, Progress in Earth and Planetary Science, 6,
1082      https://doi.org/10.1186/s40645-019-0304-z, 2019.

1083  Stoner, A. M. K., Hayhoe, K., and Wuebbles, D. J.: Assessing general circulation model simulations of atmospheric
1084      teleconnection patterns, Journal of Climate, 22, 4348–4372, https://doi.org/10.1175/2009jcli2577.1, 2009.

1085  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,
1086      Y., Lee, J., Lee, J., Kim, J.-S., Marzin, C., and Byun, Y.-H.: Climate change projection in the Twenty-First
1087      Century simulated by NIMS-KMA CMIP6 model based on new GHGs concentration pathways, Asia-pacific
1088      Journal of Atmospheric Sciences, 57, 851–862, https://doi.org/10.1007/s13143-021-00225-6, 2021.

1089  Tang, Q., Prather, M. J., Hsu, J., Ruiz, D. J., Cameron-Smith, P. J., Xie, S., and Golaz, J.-C.: Evaluation of the
1090      interactive stratospheric ozone (O3v2) module in the E3SM version 1 Earth system model, Geoscientific
1091      Model Development, 14, 1219–1236, https://doi.org/10.5194/gmd-14-1219-2021, 2021.

1092  Tang, S., Fast, J. D., Zhang, K., Hardin, J. C., Varble, A. C., Shilling, J. E., Mei, F., Zawadowicz, M. A., and Ma, P.-
1093      L.: Earth System Model Aerosol–Cloud Diagnostics (ESMAC Diags) package, version 1: assessing E3SM
1094      aerosol predictions using aircraft, ship, and surface measurements, Geosci. Model Dev., 15, 4055–4076,
1095      https://doi.org/10.5194/gmd-15-4055-2022, 2022.

1096  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.:
1097      Earth System Model Aerosol–Cloud Diagnostics (ESMAC Diags) package, version 2: assessing aerosols,
1098      clouds, and aerosol–cloud interactions via field campaign and long-term observations, Geosci. Model Dev.,
1099      16, 6355–6376, https://doi.org/10.5194/gmd-16-6355-2023, 2023.

1100  Taylor, K. E.: Truly Conserving with Conservative Remapping Methods, Geosci. Model Dev. Discuss. [preprint],
1101      https://doi.org/10.5194/gmd-2023-177, in review, 2023.

1102  Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, Journal of Geophysical
1103      Research, 106, 7183–7192, https://doi.org/10.1029/2000jd900719, 2001.

1104  Taylor, K. E., Stouffer, R. J., and Meehl, G. A.: An overview of CMIP5 and the experiment design, Bulletin of the
1105      American Meteorological Society, 93, 485–498, https://doi.org/10.1175/bams-d-11-00094.1, 2012.

1106  Teixeira, J., Waliser, D. E., Ferraro, R., Gleckler, P. J., Lee, T., and Potter, G. L.: Satellite observations for CMIP5:
1107      The Genesis of OBS4MIPs, Bulletin of the American Meteorological Society, 95, 1329–1334,
1108      https://doi.org/10.1175/bams-d-12-00204.1, 2014.

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

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

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

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

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

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

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

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

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

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

Whitehall, K., Mattmann, C., Waliser, D., Kim, J., Goodale, C., Hart, A., Ramirez, P., Zimdars, P., Crichton, D.,
Jenkins, G., Jones, C., Asrar, G., and Hewitson, B.: Building Model Evaluation and Decision Support
Capacity for CORDEX, WMO Bulletin, 61, available at:

1145    https://public.wmo.int/en/resources/bulletin/building-model-evaluation-and-decision-support-capacity-
1146        cordex (last access date: 14 September 2023), 2012.

1147 Williams, D. N.: Visualization and analysis tools for ultrascale climate data, Eos, Transactions American Geophysical
1148        Union, 95, 377–378, https://doi.org/10.1002/2014eo420002, 2014.

1149 Williams, D. N., Balaji, V., Cinquini, L., Denvil, S., Duffy, D. Q., Evans, B., Ferraro, R., Hansen, R., Lautenschlager,
1150        M., and Trenham, C.: A global repository for Planet-Sized experiments and observations, Bulletin of the
1151        American Meteorological Society, 97, 803–816, https://doi.org/10.1175/bams-d-15-00132.1, 2016.

1152 Williams, D. N., Doutriaux, C., Drach, R., and McCoy, R.: The Flexible Climate Data Analysis Tools (CDAT) for
1153        Multi-model Climate Simulation Data, IEEE International Conference on Data Mining Workshops, 254–261,
1154        https://doi.org/10.1109/icdmw.2009.64, 2009.

1155 Wong, P. C., Shen, H. W., Leung, R., Hagos, S., Lee, T. Y., Tong, X. and Lu, K.: Visual analytics of large-scale
1156        climate model data, in: 2014 IEEE 4th Symposium on Large Data Analysis and Visualization (LDAV), 85-
1157        92, https://doi.org/10.1109/LDAV.2014.7013208, 2014.

1158 Xie, P., Joyce, R., Wu, S., Yoo, S.H., Yarosh, Y., Sun, F. and Lin, R.: Reprocessed, bias-corrected CMORPH global
1159        high-resolution precipitation estimates from 1998, Journal of Hydrometeorology, 18, 1617-1641, 2017.

1160 Xue, Z. and Ullrich, P. A.: A Comprehensive Intermediate-Term Drought Evaluation System and Evaluation of
1161        Climate Data Products over the Conterminous United States, Journal of Hydrometeorology,
1162        https://doi.org/10.1175/jhm-d-20-0314.1, 2021.

1163 Young, A. H., Knapp, K. R., Inamdar, A. K., Hankins, W., and Rossow, W. B.: The International Satellite Cloud
1164        Climatology Project H-Series climate data record product, Earth System Science Data, 10, 583–593,
1165        https://doi.org/10.5194/essd-10-583-2018, 2018.

1166 Zelinka, M. D., Klein, S. A., Qin, Y., and Myers, T. A.: Evaluating climate models' cloud feedbacks against expert
1167        judgment, Journal of Geophysical Research: Atmospheres, 127, https://doi.org/10.1029/2021jd035198,
1168        2022.

1169 Zelinka, M. D., Myers, T. A., McCoy, D. T., Po-Chedley, S., Caldwell, P. M., Ceppi, P., Klein, S. A., and Taylor, K.
1170        E.: Causes of higher climate sensitivity in CMIP6 models, Geophysical Research Letters, 47,
1171        e2019GL085782, https://doi.org/10.1029/2019GL085782, 2020.

1172 Zhang, C., Golaz, J.-C., Forsyth, R., Vo, T., Xie, S., Shaheen, Z., Potter, G. L., Asay-Davis, X. S., Zender, C. S., Lin,
1173        W., Chen, C.-C., Terai, C. R., Mahajan, S., Zhou, T., Balaguru, K., Tang, Q., Tao, C., Zhang, Y.,
1174        Emmenegger, T., Burrows, S., and Ullrich, P. A.: The E3SM Diagnostics Package (E3SM Diags v2.7): a
1175        Python-based diagnostics package for Earth system model evaluation, Geosci. Model Dev., 15, 9031–9056,
1176        https://doi.org/10.5194/gmd-15-9031-2022, 2022.

1177 Zhang, C. and Hendon, H. H.: Propagating and standing components of the intraseasonal oscillation in tropical
1178        convection, Journal of the Atmospheric Sciences, 54, 741–752, https://doi.org/10.1175/1520-
1179        0469(1997)054, 1997.

1180 Zhang, C., Xie, S., Klein, S. A., Ma, H. Y., Tang, S., Van Weverberg, K., Morcrette, C. J. and Petch, J.: CAUSES:
1181        Diagnosis of the summertime warm bias in CMIP5 climate models at the ARM Southern Great Plains site.
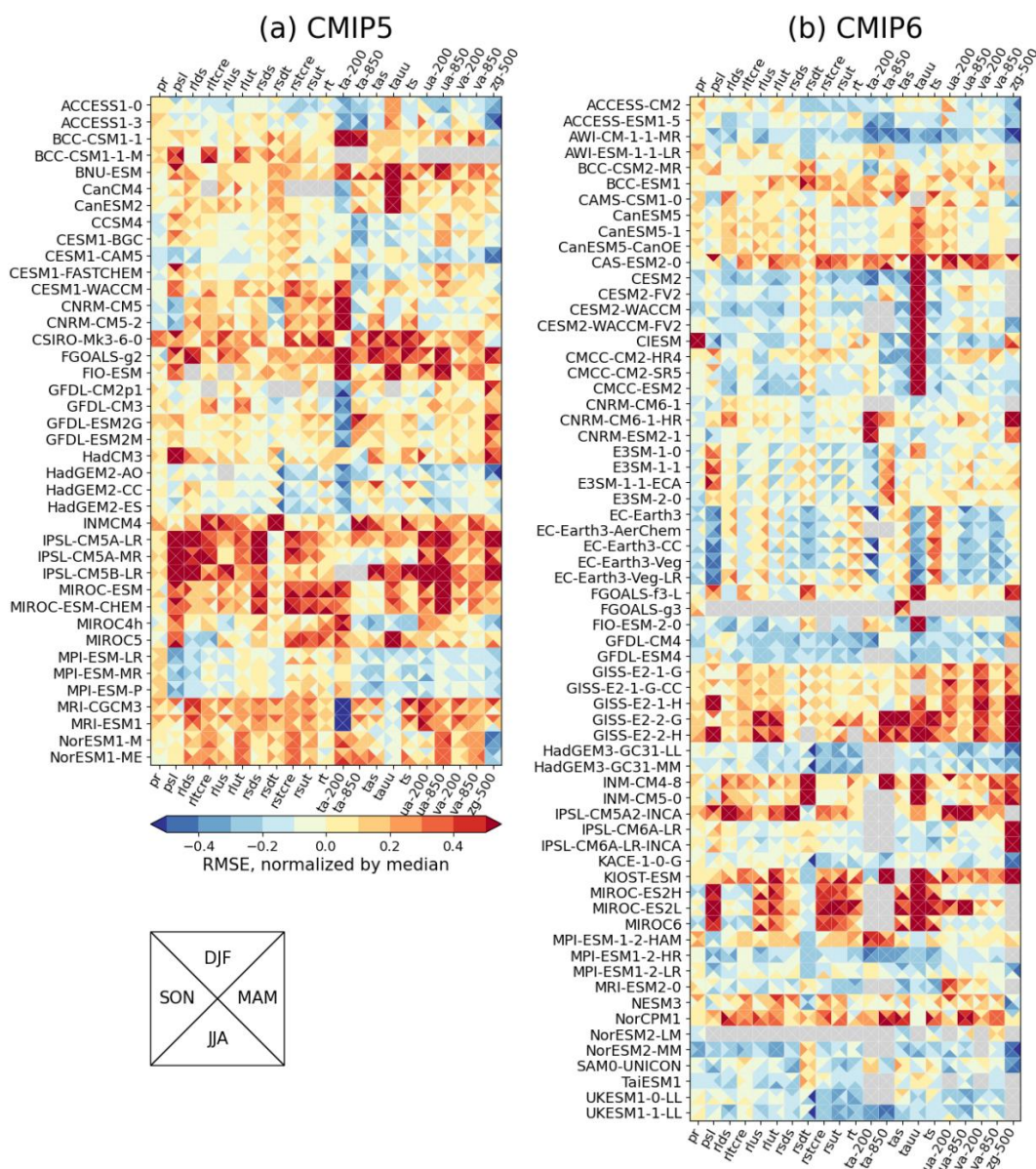
1182    Journal of Geophysical Research: Atmospheres, 123, 2968-2992, https://doi.org/10.1002/2017JD027200,
1183    2018.

1184    Zhang, C., Xie, S., Tao, C., Tang, S., Emmenegger, T., Neelin, J. D., Schiro, K. A., Lin, W. and Shaheen, Z.: The
1185    ARM data-oriented metrics and diagnostics package for climate models: A new tool for evaluating climate
1186    models with field data, https://doi.org/10.1175/BAMS-D-19-0282.1, Bulletin of the American
1187    Meteorological Society, 101, E1619-E1627, 2020.

1188    Zhang, M.-Z., Xu, Z., Han, Y., and Guo, W.: An improved multivariable integrated evaluation method and tool
1189    (MVIETool) v1.0 for multimodel intercomparison, Geoscientific Model Development, 14, 3079–3094,
1190    https://doi.org/10.5194/gmd-14-3079-2021, 2021.

1191    Zhao, B., Lin, P., Hu, A., Liu, H., Ding, M., Yu, Z., and Yu, Y.: Uncertainty in Atlantic Multidecadal Oscillation
1192    derived from different observed datasets and their possible causes, Frontiers in Marine Science, 9,
1193    https://doi.org/10.3389/fmars.2022.1007646, 2022.

1194    Zhao, M., Golaz, J.-C., Held, I. M., Guo, H., Balaji, V., Benson, R., Chen, J. H., Chen, X., Donner, L. J., Dunne, J.,
1195    Dunne, K. A., Durachta, J., Fan, S.-M., Freidenreich, S. M., Garner, S. T., Ginoux, P., Harris, L., Horowitz,
1196    L. W., Krasting, J. P., Langenhorst, A. R., Zhi, L., Lin, P., Lin, S. J., Malyshev, S., Mason, E., Milly, P. C.
1197    D., Ming, Y., Naik, V., Paulot, F., Paynter, D., Phillipps, P. J., Radhakrishnan, A., Ramaswamy, V.,
1198    Robinson, T., Schwarzkopf, D., Seman, C. J., Shevliakova, E., Shen, Z., Shin, H. H., Silvers, L. G., Wilson,
1199    J. R., Winton, M., Wittenberg, A. T., Wyman, B., and Xiang, B.: The GFDL Global Atmosphere and Land
1200    Model AM4.0/LM4.0: 1. Simulation characteristics with prescribed SSTS, Journal of Advances in Modeling
1201    Earth Systems, 10, 691–734, https://doi.org/10.1002/2017ms001208, 2018.

1202

1203 **Table 1.** List of variables and observation datasets used as reference datasets for the PMP's
1204 mean climate evaluation in this paper (Section 3.1 and Figs. 1-2). A ditto mark (") indicates the
1205 same as above.
1206

| 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 | " | |

**1208** **Figure 1**. Portrait plot for spatial RMSE (uncentered) of global seasonal climatologies for (a)
1209    CMIP5 (models ACCESS1-0 to NorESM1-ME on the ordinate) and (b) CMIP6 (models
1210    ACCESS-CM2 to UKESM1-1-LL on the ordinate) for 1981-2005 epoch. The RMSE of each
1211    variable is normalized by the median RMSE of all CMIP5 and 6 models. A result of 0.2 (-0.2) is
1212    indicative of an error that is 20% greater (lesser) than the median RMSE across all models.
1213    Models in each group are sorted in alphabetical order. Full names of variable names on the
1214    abscissa and their reference datasets can be found in Table 1. Detailed information for models
1215    can be found at the *Earth System Documentation* (ES-DOC, https://search.es-doc.org/; Pascoe

1216    et al., 2020). The interactive version of the Portrait plot in this figure is available on the PMP
1217    result pages on the PCMDI website (https://pcmdi.llnl.gov/metrics/mean_clim/).

**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/).

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

**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 rows of gray boxes. 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 whereas in the lower-right triangle, the model results are shown relative to the ERA-20C. For SST-based modes (NPGO and PDO), results are shown relative to HadISSTv1.1 (upper-left triangle) and HadISSTv2.1 (lower-right triangle). 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. The figure is adapted from Lee et al. 2021b.

1281
1282     (a) Observation



1283
1284     (b) Model



1285
1286
**Figure 5.** MJO EWR diagnostics – wavenumber-frequency power spectra – from (a) GPCP v1.3
(Huffman et al., 2001) and (b) IPSL-CM5B-LR model of CMIP5. The EWR is defined as the ratio
of eastward power (averaged in the box on the right) to westward power (averaged in the box
on the left) from the 2-dimensional wavenumber-frequency power spectra of daily 10°S–10°N
averaged precipitation in November to April (shaded, $mm^2$ $day^{-2}$). Power spectra are calculated
for each year and then averaged over all years of data. The units of power spectra for the
precipitation is $mm^2$ $day^{-2}$ per frequency interval per wavenumber interval.

**Figure 6.** MJO EWR 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.

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

1315



1316
**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
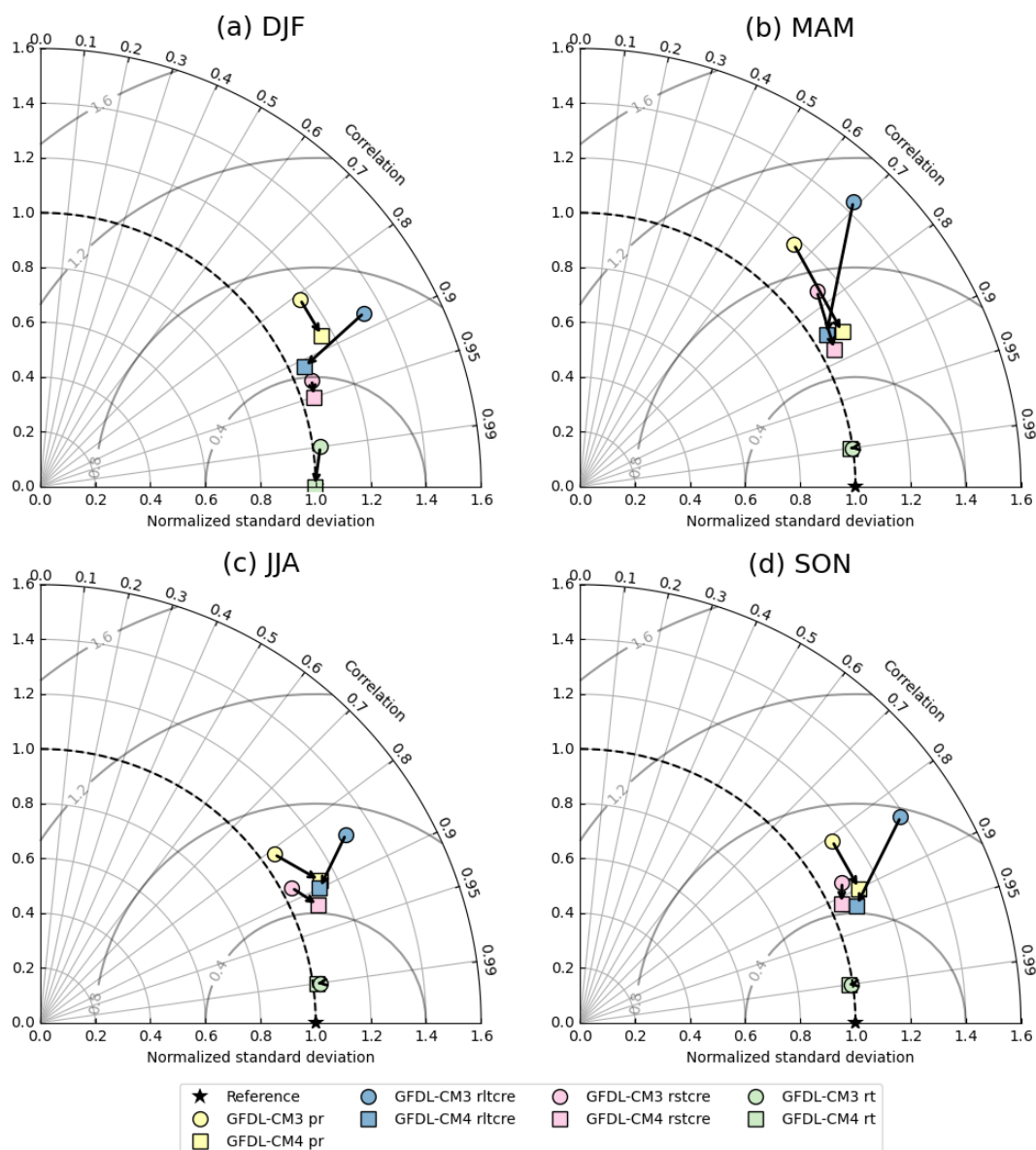illustrative model (GFDL-CM4) is highlighted.

1323



1324
1325
1326
1327 **Figure 9.** Proposed suite of baseline metrics for simulated precipitation benchmarking (figure
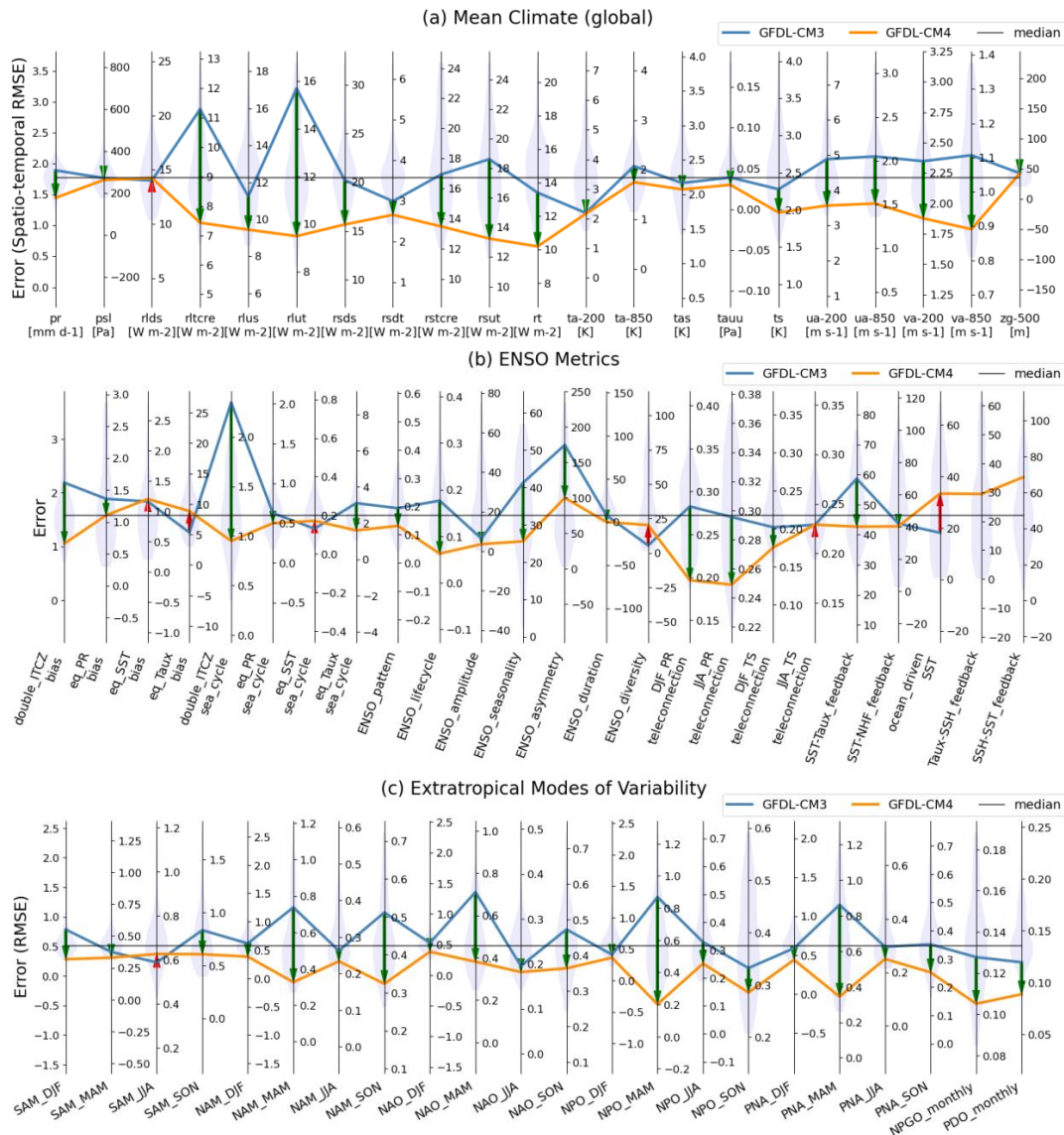1328 reprinted from workshop report; US DOE, 2020).
1329
1330

**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).

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

**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. Considering lower indicates better, 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.