

# Evaluating Precipitation Distributions at Regional Scales: A Benchmarking Framework and Application to CMIP 5 and 6 Models

Min-Seop Ahn<sup>1,2,3,\*</sup>, Paul A. Ullrich<sup>1,4</sup>, Peter J. Gleckler<sup>1</sup>, Jiwoo Lee<sup>1,\*</sup>, Ana C. Ordóñez<sup>1</sup>,  
and Angeline G. Pendergrass<sup>5,6</sup>

<sup>1</sup>*PCMDI, Lawrence Livermore National Laboratory, Livermore, CA, USA*

<sup>2</sup>*NASA Goddard Space Flight Center, Greenbelt, MD, USA*

<sup>3</sup>*ESSIC, University of Maryland, College Park, MD, USA*

<sup>4</sup>*Department of Land, Air and Water Resources, University of California, Davis, CA, USA*

<sup>5</sup>*Earth and Atmospheric Science, Cornell University, Ithaca, NY, USA*

<sup>6</sup>*National Center for Atmospheric Research, Boulder, CO, USA*

June 2023

Revised

*Geoscientific Model Development*

\* Corresponding author: Min-Seop Ahn ([ahn6@llnl.gov](mailto:ahn6@llnl.gov)) and Jiwoo Lee ([lee1043@llnl.gov](mailto:lee1043@llnl.gov))

1 **Abstract**

2 As the resolution of global Earth system models increases, regional scale evaluation is  
3 becoming ever more important. This study presents a framework for quantifying  
4 precipitation distributions at regional scales and applies it to evaluate CMIP 5 and 6  
5 models. We employ the IPCC AR6 climate reference regions over land and propose  
6 refinements to the oceanic regions based on the homogeneity of precipitation distribution  
7 characteristics. The homogeneous regions are identified as heavy, moderate, and light  
8 precipitating areas by K-means clustering of IMERG precipitation frequency and amount  
9 distributions. With the global domain partitioned into 62 regions, including 46 land and 16  
10 ocean regions, we apply 10 established precipitation distribution metrics. The collection  
11 includes metrics focused on the maximum peak, lower 10th percentile, and upper 90th  
12 percentile in precipitation amount and frequency distributions, the similarity between  
13 observed and modeled frequency distributions, an unevenness measure based on  
14 cumulative amount, average total intensity on all days with precipitation, and number of  
15 precipitating days each year. We apply our framework to 25 CMIP5 and 41 CMIP6  
16 models, and 6 observation-based products of daily precipitation. Our results indicate that  
17 many CMIP 5 and 6 models substantially overestimate the observed light precipitation  
18 amount and frequency as well as the number of precipitating days, especially over mid-  
19 latitude regions outside of some land regions in the Americas and Eurasia. Improvement  
20 from CMIP 5 to 6 is shown in some regions, especially in mid-latitude regions, but it is not  
21 evident globally, and over the tropics most metrics point toward degradation.

## 22 1. Introduction

23 Precipitation is a fundamental characteristic of the Earth's hydrological cycle and one that  
24 can have large impacts on human activity. The impact of precipitation depends on its  
25 intensity and frequency characteristics (e.g., Trenberth et al. 2003; Sun et al. 2006;  
26 Trenberth and Zhang 2018). Even with the same amount of precipitation, more intense  
27 and less frequent rainfall is more likely to lead to extreme precipitation events such as  
28 floods and drought compared to less intense and more frequent rainfall. While mean  
29 precipitation has improved in Earth system models, the precipitation distributions continue  
30 to have biases (e.g., Dai 2006; Fiedler et al. 2020), which limits the utility of these  
31 simulations, especially at the level of accuracy that is increasingly demanded in order to  
32 anticipate and adapt to changes in precipitation due to global warming.

33

34 Multi-model intercomparison with a well-established diagnosis framework facilitates  
35 identifying common model biases and sometimes yields insights into how to improve  
36 models. The Coupled Model Intercomparison Project (CMIP; Meehl et al. 2000, 2005,  
37 2007; Taylor et al. 2012; Eyring et al. 2016) is a well-established experimental protocol to  
38 intercompare state-of-the-art Earth system models, and the number of models and  
39 realizations participating in CMIP has been growing through several phases from 1  
40 (Meehl et al. 2000) to 6 (Eyring et al. 2016). Given the increasing number of models,  
41 developed at higher resolution and with increased complexity, modelers and analysts  
42 could benefit from capabilities that help synthesize the consistency between observed  
43 and simulated precipitation. As discussed in previous studies (e.g., Abramowitz 2012),  
44 our reference to model benchmarking implies model evaluation with community-

45 established reference data sets, performance tests (metrics), variables, and spatial and  
46 temporal resolutions. The U.S. Department of Energy (DOE) envisioned a framework for  
47 both baseline and exploratory precipitation benchmarks (U.S. DOE. 2020) as summarized  
48 by Pendergrass et al. (2020). While the exploratory benchmarks focus on process-  
49 oriented and phenomena-based metrics at a variety of spatiotemporal scales (Leung et  
50 al. 2022), the baseline benchmarks target well-established measures such as mean state,  
51 the seasonal and diurnal cycles, variability across timescales, intensity/frequency  
52 distributions, extremes, and drought (e.g., Gleckler et al. 2008; Covey et al. 2016; Wehner  
53 et al. 2020; Ahn et al. 2022). The current study builds on the baseline benchmarks by  
54 proposing a framework for benchmarking simulated precipitation distributions against  
55 multiple observations using well-established metrics and reference regions. To ensure  
56 consistent application of this framework, the metrics used herein are implemented and  
57 made available as part of the widely-used Program for Climate Model Diagnosis and  
58 Intercomparison (PCMDI) metrics package.

59  
60 Diagnosing precipitation distributions and formulating metrics that extract critical  
61 information from precipitation distributions have been addressed in many previous  
62 studies. Pendergrass and Deser (2017) proposed several precipitation distribution  
63 metrics based on frequency and amount distribution curves. The precipitation frequency  
64 distribution quantifies how often rain occurs at different rain rates, whereas the  
65 precipitation amount distribution quantifies how much rain falls at different rain rates.  
66 Based on the distribution curves, Pendergrass and Deser (2017) extracted rain frequency  
67 peak and amount peak where the maximum non-zero rain frequency and amount occur,

68 respectively. Pendergrass and Knutti (2018) introduced a metric that measures the  
69 unevenness of daily precipitation based on the cumulative amount curve. Their  
70 unevenness metric is defined as the number of wettest days that constitute half of the  
71 annual precipitation. In the median of station observations equatorward of 50° latitude,  
72 half of the annual precipitation falls in only about the heaviest 12 days, and generally  
73 models underestimate the observed unevenness (Pendergrass and Knutti 2018). In  
74 addition, several metrics have been developed to distill important precipitation  
75 characteristics, such as the fraction of precipitating days and simple daily intensity index  
76 (SDII, Zhang et al. 2011). In this study we implement all these well-established metrics  
77 and several other complementary metrics into our framework.

78  
79 Many studies have analyzed the precipitation distributions over large domains (e.g., Dai  
80 2006; Pendergrass and Hartmann 2014; Ma et al. 2022). Often, these domains comprise  
81 both heavily precipitating and dry regions. Given the emphasis on regional scale analysis  
82 continues to grow as models' horizontal resolution increases, interpretation of domain-  
83 averaged distributions could be simplified by defining regions that are not overly complex  
84 or heterogeneous in terms of their precipitation distribution characteristics. Iturbide et al.  
85 (2020) has identified climate reference regions that have been adopted in the sixth  
86 assessment report (AR6) of the Intergovernmental Panel on Climate Change (IPCC). Our  
87 framework is based on these IPCC AR6 reference regions for objective examination of  
88 precipitation distributions over land. Over the ocean we have revised some of the regions  
89 of Iturbide et al. (2020) to better isolate homogeneous precipitation distribution  
90 characteristics.

91

92 In this study, we propose a modified IPCC AR6 reference regions and a framework for  
93 regional scale quantification of simulated precipitation distributions, which is implemented  
94 into the PCMDI metrics package to enable researchers to readily use the metric collection  
95 in a common framework. The remainder of the paper is organized as follows: Sections 2  
96 and 3 describe the data and analysis methods. Section 4 presents results including the  
97 application and modification of IPCC AR6 climate reference regions, evaluation of CMIP  
98 5 and 6 models with multiple observations, and their improvement across generations.  
99 Sections 5 and 6 discuss and summarize the main accomplishments and findings from  
100 this study.

101

102

## 103 **2. Data**

### 104 2.1. Observational data

105 For reference data, we use six global daily precipitation products first made available as  
106 part of the Frequent Rainfall Observations on GridS (FROGS) database (Roca et al.,  
107 2019) and then further aligned with CMIP output via the data specifications of the  
108 Observations for Model Intercomparison Project (Obs4MIPs, Waliser et al. 2020). These  
109 include five satellite-based products and a recent atmospheric reanalysis product. The  
110 satellite-based precipitation products include the Integrated Multi-satellitE Retrievals for  
111 GPM version 6 final run product (Huffman et al. 2020; hereafter IMERG), the Tropical  
112 Rainfall Measuring Mission Multi-satellite Precipitation Analysis 3B42 version 7 product  
113 (Huffman et al. 2007; hereafter TRMM), the bias-corrected Climate Prediction Center

114 Morphing technique product (Xie et al. 2017; hereafter CMORPH), the Global  
115 Precipitation Climatology Project 1DD version 1.3 (Huffman et al. 2001; hereafter GPCP),  
116 and Precipitation Estimation from Remotely Sensed Information using Artificial Neural  
117 Networks (Ashouri et al. 2015; hereafter PERSIANN). The reanalysis product included  
118 for context is the European Centre for Medium-Range Weather Forecasts (ECMWF)'s  
119 fifth generation of atmospheric reanalysis (Hersbach et al. 2020; hereafter ERA5). Table  
120 1 summarizes the observational datasets with the data source, coverage of domain and  
121 period, resolution of horizontal space and time frequency, and references. We use the  
122 data periods available via FROGS and Obs4MIPs as follows: 2001-2020 for IMERG,  
123 1998-2018 for TRMM, 1998-2012 for CMORPH, 1997-2020 for GPCP, 1984-2018 for  
124 PERSIANN, and 1979-2018 for ERA5.

125

## 126 2.2. CMIP model simulations

127 We analyze daily precipitation from all realizations of AMIP simulations available from  
128 CMIP5 (Taylor et al. 2012) and CMIP6 (Eyring et al. 2016). We have chosen to  
129 concentrate our analysis on AMIP simulations rather than the coupled Historical  
130 simulations because the simulated precipitation in the latter is strongly influenced by  
131 biases in the modeled sea surface temperature, complicating any interpretation regarding  
132 the underlying causes of the precipitation errors. Table 2 lists the participating models,  
133 the number of realizations, and the horizontal resolution in each modeling institute. We  
134 evaluate the most recent 20 years (1985-2004) that both CMIP 5 and 6 models have in  
135 common for a fair comparison with satellite-based observations.

136

137

### 138 **3. Methods**

139 In our framework we apply 10 metrics that characterize different and complementary  
140 aspects of the intensity distribution of precipitation at regional scales. Table 3 summarizes  
141 the metrics including their definition, purpose, and references. The computation of the  
142 metrics has been implemented and applied in the PCMDI metrics package (PMP;  
143 Gleckler et al. 2008, 2016).

144

#### 145 3.1. Frequency and amount distributions

146 Following Pendergrass and Hartmann (2014) and Pendergrass and Deser (2017), we use  
147 logarithmically-spaced bins of daily precipitation to calculate both the precipitation  
148 frequency and amount distributions. Each bin is 7% wider than the previous one, and the  
149 smallest non-zero bin is centered at 0.03 mm/day. The frequency distribution is the  
150 number of days in each bin normalized by the total number of days, and the amount  
151 distribution is the sum of accumulated precipitation in each bin normalized by the total  
152 number of days. Based on these distributions (Fig. 1a), we identify the rain rate where the  
153 maximum peak of the distribution appears (Amount/Frequency Peak, Pendergrass and  
154 Deser 2017; also called mode, Kooperman et al., 2016) and formulate several  
155 complementary metrics that measure the fraction of the distribution lower 10 percentile  
156 (P10) and upper 90 percentile (P90). The precipitation bins less than 0.1 mm/day are  
157 considered dry for the purpose of these calculations. The threshold rain rates for 10th and  
158 90th percentiles are defined from the amount distribution as determined from  
159 observations. Here we use IMERG as the default reference observational dataset. The



160 final frequency related metric we employ is the Perkins score, which measures the  
161 similarity between observed and modeled frequency distributions (Perkins et al. 2007).  
162 With the sum of a frequency distribution across all bins being unity, the Perkins score is  
163 defined as the sum of minimum values between observed and modeled frequency across  
164 all bins:  $Perkins\ Score = \sum_1^n \text{minimum}(Z_o, Z_m)$  where  $n$  is the number of bins,  $Z_o$  and  $Z_m$   
165 are the frequency in a given bin for observation and model, respectively. The Perkins  
166 score is a unitless scalar varying from 0 (low similarity) to 1 (high similarity).

167

### 168 3.2. Cumulative fraction of annual precipitation amount

169 Following Pendergrass and Knutti (2018), we calculate the cumulative sum of daily  
170 precipitation each year sorted in descending order (i.e., wettest to driest) and normalized  
171 by the total precipitation for that year. From the distribution for each individual year (see  
172 Fig. 1b), we obtain the metrics gauging the number of wettest days for half of annual  
173 precipitation (Unevenness, Pendergrass and Knutti 2018) and the fraction of the number  
174 of precipitating ( $\geq 1$  mm/day) days (FracPRdays). To facilitate comparison against longer-  
175 established analyses (e.g., ETCCDI, Zhang et al., 2011), we include the daily  
176 precipitation intensity, calculated by dividing the annual total precipitation by the number  
177 of precipitating days (SDII, Zhang et al. 2011). To obtain values of these metrics over  
178 multiple years, we take the median across years following Pendergrass and Knutti (2018;  
179 for unevenness).

180

### 181 3.3. Reference regions

182 We use the spatial homogeneity of precipitation characteristics as a basis for defining  
183 regions, as in previous studies (e.g., Swenson and Grotjahn 2019). In addition to  
184 providing more physically-based results, this also simplifies interpretation with robust  
185 diagnostics when we average a distribution characteristic across the region. We use K-  
186 means clustering (MacQueen 1967) with the concatenated frequency and amount  
187 distributions of IMERG over the global domain to identify homogeneous regions for  
188 precipitation distributions. K-means clustering is an unsupervised machine learning  
189 algorithm that separates characteristics of a dataset into a given number of clusters  
190 without explicitly provided criteria. This method has been widely used because it is faster  
191 and simpler than other methods. Here, we use 3 clusters to define heavy, moderate and  
192 light precipitation regions. Figure 2 shows the spatial pattern of IMERG precipitation mean  
193 state and clustering results defining heavy (blue), moderate (green), and light (orange)  
194 precipitation regions. The spatial pattern of these clustering regions resembles the pattern  
195 of the mean state of precipitation, providing a sanity check indicating that the cluster-  
196 based regions are physically reasonable. Note that the clustering result with frequency  
197 and amount distributions is not significantly altered if we incorporate cumulative amount  
198 fraction. However, the inclusion of the cumulative amount fraction to the clustering yields  
199 a slightly noisier pattern, and thus we have chosen to use the clustering result only with  
200 frequency and amount distributions.

201

202 In support of the AR6, the IPCC proposed a set of climate reference regions (Iturbide et  
203 al. 2020). These regions were defined based on geographical and political boundaries  
204 and the climatic consistency of temperature and precipitation in current climate and

205 climate change projections. When defining regions, the land regions use both information  
206 from current climate and climate change projections, while the ocean regions use only  
207 the information from climate change projections. In other words, the climatic consistency  
208 of precipitation in the current climate is not explicitly represented in the definition of the  
209 oceanic regions. Figure 3a shows the IPCC AR6 climate reference regions superimposed  
210 on our precipitation clustering regions shown in Fig. 2b. The land regions correspond  
211 reasonably well to the clustering regions, but some ocean regions are too broad, including  
212 both heavy and light precipitating regions (Fig. 3a). In this study, the ocean regions are  
213 modified based on the clustering regions, while the land regions remain the same as in  
214 the AR6 (Fig. 3b).

215

216 In the Pacific Ocean region, the original IPCC AR6 regions consist of equatorial Pacific  
217 Ocean (EPO), northern Pacific Ocean (NPO), and southern Pacific Ocean (SPO). Each  
218 of these regions includes areas of both heavy and light precipitation. EPO includes the  
219 Intertropical Convergence Zone (ITCZ), the South Pacific Convergence Zone (SPCZ),  
220 and also the dry southeast Pacific region. The NPO region includes the north Pacific storm  
221 track and the dry northeast Pacific. The SPO region includes the southern part of SPCZ  
222 and the dry southeast area of the Pacific. In our modified IPCC AR6 regions, the Pacific  
223 Ocean region is divided into four heavy precipitating regions (NPO, NWPO, PITCZ, and  
224 SWPO) and two light and moderate precipitating regions (NEPO and SEPO). Similarly,  
225 in the Atlantic Ocean region, the original IPCC AR6 regions consist of the equatorial  
226 Atlantic Ocean (EAO), northern Atlantic Ocean (NAO), and southern Atlantic Ocean  
227 (SAO), with each including both heavy and light precipitating regions. Our modified

228 Atlantic Ocean region consists of two heavy precipitating regions (NAO and AITCZ) and  
229 two light and moderate precipitating regions (NEAO and SAO). The Indian Ocean (IO)  
230 region is not modified as the original IPCC AR6 climate reference region separates well  
231 the heavy precipitating equatorial IO (EIO) region from the moderate and light  
232 precipitating southern IO (SIO) region. The Southern Ocean (SOO) is modified to mainly  
233 include the heavy precipitation region around the Antarctic. The original IPCC AR6  
234 climate reference regions consist of 58 regions including 12 oceanic regions and 46 land  
235 regions, while our modification consists of 62 regions including 16 oceanic regions and  
236 the same land regions as the original (see Table 4). Note that the Caribbean (CAR), the  
237 Mediterranean (MED), and Southeast Asia (SEA) are not counted for the oceanic regions.

238

#### 239 3.4. Evaluating model performance

240 We use two simple measures to compare the collection of CMIP 5 and 6 model  
241 simulations with the five satellite-based observational products (IMERG, TRMM,  
242 CMORPH, GPCP, and PERSIANN). One gauges how many models within the multi-  
243 model ensemble fall within the observational range between the minimum and maximum  
244 observed values for each metric and each region. Another is how many models  
245 underestimate or overestimate all observations, i.e., are outside the bounds spanned by  
246 the minimum and maximum values across the five satellite-based products. To quantify  
247 the dominance of underestimation versus overestimation of the multi-model ensemble  
248 with a single number, we use the following measure formulation:  $(nO - nU)/nT$  where  $nO$   
249 is the number of overestimating models,  $nU$  is the number of underestimating models,  
250 and  $nT$  is the total number of models. Thus, positive values represent overestimation, and

251 negative values represent underestimation. If models are mostly within the observational  
252 range or widely distributed from underestimation to overestimation, the quantification  
253 value would approach zero.

254

255 Many metrics that can be used to characterize precipitation, including those used here,  
256 are sensitive to the spatial and temporal resolutions at which the model and observational  
257 data are analyzed (e.g., Pendergrass and Knutti 2018, Chen and Dai 2019). As in many  
258 previous studies the diagnosis of precipitation in CMIP 5 and 6 models (e.g., Fiedler et al.  
259 2020; Tang et al. 2021; Ahn et al. 2022), to ensure appropriate comparisons, we conduct  
260 all analyses at a common horizontal grid of 2x2 degrees with a conservative regridding  
261 method. For models with multiple ensemble members, we first compute the metrics for  
262 all available realizations and then average the results across the realizations.

263

264

## 265 **4. Results**

### 266 4.1. Homogeneity within reference regions

267 For the regional scale analysis, we employ the IPCC AR6 climate reference regions  
268 (Iturbide et al. 2020) while we revise the region dividings over the oceans based on  
269 clustered precipitation characteristics as described in section 3.3. To quantitatively  
270 evaluate the homogeneity of precipitating distributions in the reference regions, we use  
271 three homogeneity metrics: the Perkins score (Perkins et al. 2007), Kolmogorov–Smirnov  
272 test (K-S test, Chakravart et al. 1967), and Anderson-Darling test (A-D test, Stephens  
273 1974). The three metrics measure the similarity between the regionally-averaged and

274 individual grid cell frequency distributions within the region. The Perkins score measures  
275 the overall similarity between two frequency distributions, which is one of our distribution  
276 performance metrics described in Section 3.1. The K-S and A-D tests focus more on the  
277 similarity in the center and the side of the frequency distribution, respectively. The three  
278 homogeneity metrics could complement each other as their main focuses are on different  
279 aspects of frequency distributions.

280

281 In the original IPCC AR6 reference regions, the oceanic regions show relatively low  
282 homogeneity of precipitating characteristics compared to land regions (Fig. 4). The Pacific  
283 and Atlantic Ocean regions show much lower homogeneity than the Indian Ocean,  
284 especially in EPO and EAO regions. In the modified oceanic regions, the homogeneities  
285 show an overall improvement with the three homogeneity metrics. In particular, the  
286 homogeneity over the heavy precipitating regions where the homogeneity was lower (e.g.,  
287 Pacific and Atlantic ITCZ and mid-latitude storm track regions) are largely improved. The  
288 clustering regions shown here are obtained based on IMERG precipitation. However,  
289 since different satellite-based products show substantial discrepancies in precipitation  
290 distributions, it is important to assess whether the improved homogeneity in the modified  
291 regions is similarly improved across other different datasets. Figure 5 shows the  
292 homogeneity of precipitation distribution characteristics for different observational  
293 datasets using the Perkins score. Although the region modifications we have made are  
294 based on the clustering regions of IMERG precipitation, Fig. 5 suggests that the  
295 improvement of the homogeneity over the modified regions is consistent across different  
296 observational datasets. We further tested the homogeneity for different seasons (see Fig.

297 S1 in the supplement material). The homogeneity is overall improved in the modified  
298 regions across the seasons even though we defined the reference regions based on  
299 annual data.

300

#### 301 4.2. Regional evaluation of model simulations against multiple observations

302 The precipitation distribution metrics used in this study are mainly calculated from three  
303 curves: amount distribution, frequency distribution, and cumulative amount fraction  
304 curves. Figure 6 shows these curves for three selected regions based on the clustered  
305 precipitating characteristics (NWPO, which is a heavy precipitation dominated ocean  
306 region; SEPO, a light precipitation dominated ocean region; and ENA, a heavy  
307 precipitation dominated land region). The heavy and light precipitating regions are well  
308 distinguished by their overlaid distribution curves. The amount distribution has a  
309 distinctive peak in the heavy precipitating region (Figs. 6a and 6g), while it is flatter in the  
310 light precipitating region (Fig. 6d). The frequency distribution is more centered on the  
311 heavier precipitation side in the heavy precipitating region (Figs. 6b, 6h) than in the light  
312 precipitating region (Fig 6e). The cumulative fraction increases more steeply in the light  
313 precipitating region (Fig. 6f) than in the heavy precipitating region (Figs. 6c and 6i),  
314 indicating there are fewer precipitating days in the light precipitating region. NWPO and  
315 SEPO were commonly averaged for representing the tropical ocean region in many  
316 studies, but these different characteristics in the precipitation distributions demonstrate  
317 the additional information available via a regional scale analysis. Although satellite-based  
318 observations are less reliable over the light precipitating ocean regions (e.g., SEPO), the  
319 differences between heavy and light precipitation regions are well distinguishable.

320

321 In the precipitation frequency distribution, many models show a bimodal distribution in the  
322 heavy precipitating tropical ocean region (Fig. 6b) but not in the light precipitating  
323 subtropical ocean region (Fig. 6e) or the heavy precipitating mid-latitude land region (Fig.  
324 6h). The bimodal frequency distribution is a commonly found in models and is seemingly  
325 independent of resolution (e.g., Lin et al. 2013; Kooperman et al. 2018; Chen et al. 2021;  
326 Ma et al. 2022; Martinez-Villalobos et al. 2022; Ahn et al. 2023). Ma et al. (2022)  
327 compared the frequency distributions in AMIP and HighResMIP (High Resolution Model  
328 Intercomparison Project, Haarsma et al. 2016) from CMIP6 and DYAMOND (DYnamics  
329 of the Atmospheric general circulation Modeled On Non-hydrostatic Domains, Satoh et  
330 al. 2019; Stevens et al. 2019) models, where they showed that the bimodal frequency  
331 distribution appears in many AMIP (~100km), HighResMIP (~50km), and even  
332 DYAMOND (~4km) models. Ahn et al. (2023) further compared between DYAMOND  
333 model simulations with and without a convective parameterization and showed that most  
334 DYAMOND model simulations exhibiting the bimodal distribution use a convective  
335 parameterization. ERA5 reanalysis also shows a bimodal frequency distribution (Fig. 6b),  
336 which is not surprising considering that the reproduced precipitation in ERA5 heavily  
337 depends on the model, thus exhibits this common model behavior. Because of the heavy  
338 reliance on model physics to generate its precipitation (as opposed to fields like wind, for  
339 which observations are directly assimilated), in this study we do not include ERA5  
340 precipitation among the observational products used for model evaluation.

341



342 Based on the precipitation amount, frequency, and cumulative amount fraction curves,  
343 we calculate 10 metrics (Amount peak, Amount P10, Amount P90, Frequency peak,  
344 Frequency P10, Frequency P90, Unevenness, FracPRdays, SDII, and Perkins score) as  
345 described in Section 3. Figure 7 shows the metrics with the modified IPCC AR6 climate  
346 reference regions for satellite-based observations (black), ERA5 (gray), CMIP5 (blue),  
347 and CMIP6 (red) models. The metric values vary widely across regions, especially in  
348 Amount peak, Frequency peak, Unevenness, FracPRdays, and SDII, demonstrating the  
349 additional detail provided by regional-scale precipitation-distribution metrics. In terms of  
350 the metrics based on the amount distribution (Fig. 7a-c), many models tend to simulate  
351 an Amount peak that is too light, an Amount P10 that is too high, and an Amount P90 that  
352 is too low compared to the observations, moreso in oceanic regions (regions 47-62) than  
353 in land regions. Similarly for the metrics based on the frequency distribution (Fig. 7d-f),  
354 many models show light Frequency peaks, overestimated Frequency P10, and  
355 underestimated Frequency P90 compared to observations. The similarity between  
356 frequency distribution curves (i.e., Perkins score) is higher in land regions than in ocean  
357 regions. Also, many models overestimate Unevenness and FracPRdays and  
358 underestimate SDII. These results indicate that overall, models simulate more frequent  
359 weak precipitation and less heavy precipitation compared to the observations, consistent  
360 with many previous studies (e.g., Dai 2006; Pendergrass and Hartmann 2014; Trenberth  
361 et al. 2017; Chen et al. 2021; Ma et al. 2022).

362

363 As expected from previous work, observations disagree substantially in some regions  
364 (e.g., polar and high latitude regions) and/or for some metrics (e.g., Amount P90,

365 Frequency P90). In some cases the observational spread is much larger than that of the  
366 models. We examine the observational discrepancy or spread by the ratio between the  
367 standard deviation of the five satellite-based observations (IMERG, TRMM, CMORPH,  
368 GPCP, PERSIANN) and the standard deviation of all CMIP 5 and 6 models (Fig. 8). The  
369 standard deviation of observations is much larger near polar regions and high latitude  
370 regions compared to the models' standard deviation for most metrics, as expected from  
371 the orbital configurations of the most relevant satellite constellations for precipitation  
372 (which exclude high latitudes). The Amount P90 and Frequency P90 metrics show the  
373 largest observational discrepancy among the metrics, with standard deviations of 1.5 to  
374 3 times larger over some high latitude regions and about 3-8 times larger over polar  
375 regions in observations compared to the models. On the other hand, Unevenness,  
376 FracPRdays, and Amount P10 show the least observational discrepancy – the models'  
377 standard deviation is about 2-8 times larger than for observations over some tropical and  
378 subtropical regions; nonetheless, the standard deviation among observations is larger  
379 over most of the high latitude and polar regions. Model evaluation in the regions with large  
380 disagreement among observational products remains a challenge. Note that the standard  
381 deviation of five observations would be sensitive as there are outlier observations for  
382 some regions and metrics (e.g., many ocean regions in Amount P90). Moreover,  
383 observational uncertainties are rarely well quantified or understood, so agreements  
384 among observational datasets may not always allow us to rule out common errors among  
385 observations (e.g., for warm light precipitation over the subtropical ocean).

386

387 To attempt to account for discrepancies among observational datasets in the model  
388 evaluation framework, we use two different approaches to evaluate model performance  
389 with multiple observations, as described in Section 3.4. The first approach is to assess  
390 the number of models that are within the observational range. Figure 9 shows the CMIP6  
391 model evaluation with each metric, and the regions where the standard deviation among  
392 observations is larger than among models are stippled gray to avoid them from the model  
393 performance evaluation. In Amount peak, some subtropical regions (e.g., ARP, EAS,  
394 NEPO, CAU, and WSAF) show relatively good model performance (more than 70% of  
395 models fall in the observational range), while some tropical and subtropical (e.g., PITCZ,  
396 AITCZ, and SEPO) and polar (e.g., RAR, EAN, and WAN) regions show poor model  
397 performance (less than 30% of models fall in observational range). For Amount P10,  
398 many regions are poorly captured by the simulations, except for some subtropical land  
399 regions (e.g., EAS, NCA, CAU, and WSAF). In Amount P90, most regions are uncertain  
400 (i.e., the standard deviation among observations is larger than among models) making it  
401 difficult to evaluate model performance, while some tropical regions near the Indo-Pacific  
402 warmpool (EIO, SEA, NWPO, and NAU) exhibit very good model performance (more than  
403 90% of models fall in observational range). In the Frequency metrics (peak, P10, and  
404 P90), more regions are difficult to evaluate model performance than in Amount metrics,  
405 while in some tropical and subtropical regions (e.g., PITCZ, SWPO, NWPO, SEA, SAO,  
406 and NES) model performance is good. However, good model performance could  
407 alternatively arise from a large observational range (see Fig. 7). Unevenness,  
408 FracPRdays, SDII, and Perkins score have a smaller fraction of models within the  
409 observational range in tropical regions than the Amount and Frequency metrics. In

410 particular, fewer than 10% of CMIP6 models fall within the observational range for  
411 Unevenness and FracPRdays over some tropical oceanic regions (e.g., PITCZ, NEPO,  
412 SEPO, AITCZ, NEAO, SAO, and SIO).

413

414 Examining the fraction of CMIP5 models falling within the range of observations, CMIP5  
415 models have a spatial pattern of model performance similar to that of CMIP6 models (see  
416 Fig. S2 in supplement), and the improvement from CMIP5 to CMIP6 seems subtle. We  
417 quantitatively assess the improvement from CMIP5 to CMIP6 by subtracting the  
418 percentage of CMIP5 from CMIP6 models falling within the range of observations (Fig.  
419 10). For some metrics (e.g., Amount peak, Amount and Frequency P10, and Amount and  
420 Frequency P90) and for some tropical and subtropical regions (e.g., SEA, EAS, SAS,  
421 ARP, and SAH), improvement is apparent. Compared to CMIP5, 5-25% more CMIP6  
422 models fall in the observational range in these regions. However, for the other metrics  
423 (e.g., Frequency peak, FracPRdays, SDII, Perkins score), CMIP6 models perform  
424 somewhat worse. Over some tropical and subtropical oceanic regions (e.g., PITCZ,  
425 NEPO, AITCZ, and NEAO), 5-25% more CMIP6 than CMIP5 models are out of the  
426 observational range. This result is from all available CMIP5 and CMIP6 models, so it may  
427 reflect the fact that some models are participated in only CMIP5 or CMIP6, but not both  
428 (see Table 2). To isolate improvements that may have occurred between successive  
429 generations of the same models, we also compared only the models that participated in  
430 both CMIP5 and CMIP6 (see Fig. S3). Overall, the spatial characteristics of the  
431 improvement/degradation in CMIP6 from CMIP5 is consistent, while more degradation is

432 apparent when we compare this subset of models, especially over the tropical oceanic  
433 regions (e.g., PITCZ, AITCZ, NWPO, and SEPO).

434

435 The second approach to account for discrepancies among observations in model  
436 performance evaluation is to count the number of models that are lower or higher than all  
437 satellite-based observations for each metric and each region. Figure 11 shows the spatial  
438 patterns of the model performance evaluation with each metric for CMIP6 models.  
439 Underestimation is indicated by a negative sign, while overestimation is indicated by a  
440 positive sign via the formulation described in Section 3.4. Amount peak is overall  
441 underestimated in most regions, indicating the amount distributions in most CMIP6  
442 models are shifted to lighter precipitation compared to observations. In many regions,  
443 more than 50% of the CMIP6 models underestimate Amount peak. In particular, over  
444 many tropical and southern hemisphere ocean regions (e.g., PITCZ, AITCZ, EIO, SEPO,  
445 SAO, and SOO), more than 70% of the models underestimate the Amount peak. The  
446 underestimation of Amount peak is accompanied by overestimation of Amount P10 and  
447 underestimation of Amount P90. More than 70% of CMIP6 models overestimate Amount  
448 P10 in many oceanic regions; especially in the southern and northern Pacific and Atlantic,  
449 the southern Indian Ocean, and Southern Ocean more than 90% of the models  
450 overestimate the observed Amount P10. For Amount P90, it appears that many models  
451 fall within the observational range; however, observational range in Amount P90 (green  
452 boxes in Fig. 7c) is large and driven primarily by just one observational dataset (IMERG),  
453 especially in ocean regions.

454

455 For the frequency-based metrics (i.e., peak, P10, and P90; Figs. 11d-f), CMIP6 models  
456 show similar bias characteristics to Amount metrics (Figs. 11a-c), although performance  
457 is better than for Amount metrics. Over some tropical (e.g., NWPO, PITCZ, and SWPO )  
458 and Eurasia (e.g., EEU, WSB, and ESB) regions, less than 10% of models fall outside of  
459 the observed range. Unevenness and FracPRdays are severely overestimated in models.  
460 More than 90% of models overestimate the observed Unevenness (Fig. 11g) and  
461 FracPRdays (Fig. 11h) globally, especially over oceanic regions, consistent with  
462 Pendergrass and Knutti (2018). SDII is underestimated in many regions globally,  
463 especially in some heavily-precipitating regions (e.g., PITCZ, AITCZ, EIO, SEA, NPO,  
464 NAO, SWPO, and SOO). For the Perkins score, model simulations have poorer  
465 performance in the tropics than in the mid-latitudes and polar regions. Performance by  
466 these various metrics is generally consistent with the often-blamed too-frequent light  
467 precipitation and too rare heavy precipitation in simulations.

468  
469 The characteristics of CMIP5 compared to CMIP6 simulations (Fig. S4) show little  
470 indication of improvement. Here we quantitatively evaluate the improvement in CMIP6  
471 from CMIP5 for each metric and region. Figure 12 shows the difference between CMIP5  
472 and CMIP6 in terms of the percentage of models that under- or over-estimate each metric.  
473 In mid-latitudes, there appears to have been an improvement in performance, however in  
474 the tropics, there appears to be more degradation. Over some heavily-precipitating  
475 tropical regions (e.g., PITCZ, AITCZ, EIO, and NWPO), 10-25% more models in CMIP6  
476 than in CMIP5 overestimate Amount P10, Unevenness, and FracPRdays and  
477 underestimate/underperform on Amount peak, SDII, and Perkins score. This indicates

478 that CMIP6 models simulate more frequent light precipitation and less frequent heavy  
479 precipitation over the heavily-precipitating tropical regions. Over some mid-latitude land  
480 regions (e.g., EAS, ESB, RFE, and ENA), on the other hand, 5-20% more models in  
481 CMIP6 than in CMIP5 simulate precipitation distributions close to observations (i.e., less  
482 light precipitation and more heavy precipitation). To evaluate the improvement between  
483 model generation, we also compare only the models that participated in both CMIP5 and  
484 CMIP6 (Fig. S5) rather than all available CMIP5 and CMIP6 models. For the subset of  
485 models participating in both generations, the improvement characteristics are similar for  
486 all models, although more degradation is exhibited over some tropical oceanic regions  
487 (e.g., PITCZ, NWPO, and SWPO). This also indicates that some models newly  
488 participating in CMIP6, and not in the CMIP5, have higher than average performance.

489

#### 490 4.3. Correlation between metrics

491 Each precipitation distribution metric implemented in this study is chosen to target  
492 different aspects of the distribution of precipitation. To the extent that precipitation  
493 probability distributions are governed by a small number of key parameters (as argued by  
494 Martinez-Villalobos and Neelin 2019), we should expect additional metrics to be highly  
495 correlated. Figure 13 shows the global weighted average of correlation coefficients  
496 between the precipitation distribution metrics across CMIP5 and CMIP6 models. Higher  
497 correlation coefficients are found to be between Amount P90 and Frequency P90 (0.98)  
498 and between Amount P10 and Frequency P10 (0.67). This is expected because the  
499 amount and frequency distributions differ only by a factor of the precipitation rate (e.g.,  
500 Pendergrass and Hartmann 2014). Another higher correlation coefficient is between

501 Unevenness and FracPRdays (0.77), indicating that the number of the heaviest  
502 precipitating days for half of annual precipitation and the total number of annual  
503 precipitating days are related. Amount and Frequency peak metrics are negatively  
504 correlated to P10 metrics and positively correlated to P90 metrics, but the correlation  
505 coefficients are not very high (lower than 0.62). This is because the peak metrics focus  
506 on typical precipitation, rather than the light and heavy ends of the distribution that are  
507 the focus of P10 and P90 metrics. SDII is more negatively correlated with Amount P10 (-  
508 0.67) and positively correlated with Amount peak (0.61) and less so with Amount P90  
509 (0.48), implying that SDII is mainly influenced by weak precipitation amounts rather than  
510 heavy precipitation amounts. The Perkins score shows relatively high negative correlation  
511 with Unevenness (-0.62), FracPRdays (-0.59), and Amount P10 (-0.59). This indicates  
512 that the discrepancy between the observed and modeled frequency distributions is partly  
513 associated with the overestimated light precipitation in models. The correlation  
514 coefficients between the metrics other than those discussed above are lower than 0.6.  
515 While there is some redundant information within the collection of metrics included in our  
516 framework, we retain all metrics so that others can select an appropriate subset for their  
517 own application. This also preserves the ability to readily identify outlier behavior of an  
518 individual model across a wide range of regions and statistics.

519

#### 520 4.4. Influence of spatial resolution on metrics

521 Many metrics for the precipitation distribution are sensitive to the spatial resolution of  
522 the underlying data (e.g., Pendergrass and Knutti 2018; Chen and Dai 2019). Figure 14  
523 shows how our results (which are all based on data at 2° resolution) are impacted if we



524 calculate the metrics from data coarsened to 4° grid instead. As expected, there is clearly  
525 some sensitivity to the spatial scale at which our precipitation distribution metrics are  
526 computed, but the correlation among datasets (both models and observations) between  
527 the two resolutions is very high, indicating that evaluations at either resolution should be  
528 consistent. At the coarser resolution, Amount peak and SDII are consistently smaller (as  
529 expected); Amount P10 and Frequency P10 tend to be smaller as well. Meanwhile,  
530 Unevenness and FracPRdays are consistently large (as expected); Amount P90,  
531 Frequency P90, and Perkins score are generally larger as well. Chen and Dai (2019)  
532 discussed a grid aggregation effect that is associated with the increased probability of  
533 precipitation as the horizontal resolution becomes coarser. This effect is clearly evident  
534 with increased Unevenness (Fig. 14g), FracPRdays (Fig. 14h), and decreased SDII (Fig.  
535 14i) in coarser resolution. However, despite these differences, the relative model  
536 performance is not very sensitive to the spatial scale at which we apply our analysis. The  
537 correlation coefficients between results based on all data interpolated to 2° or 4°  
538 horizontal resolutions are above 0.9 for all of our distribution metrics. Conclusions on  
539 model performance are relatively insensitive to the target resolution.

540

541

## 542 **5. Discussion**

543 Analyzing the distribution of precipitation intensity lags behind temperature and even  
544 mean precipitation. Challenges include choosing appropriate metrics and analysis  
545 resolution to characterize this highly non-gaussian variable and interpreting model skills  
546 in the face of substantial observational uncertainty. Comparing results derived at 2° and

547 4° horizontal resolution for CMIP class models, we find that the quantitative changes in  
548 assessed performance are highly consistent across models and consequently have little  
549 impact on our conclusions. More work is needed to determine how suitable this collection  
550 of metrics may be for evaluating models with substantially higher resolutions (e.g.,  
551 HighResMIP, Haarsma et al. 2016). We note that more complex measures have been  
552 designed to be scale independent (e.g., Martinez-Villalobos and Neelin 2019; Martinez-  
553 Villalobos et al. 2022), and these may become increasingly important with continued  
554 interest in models developed at substantially higher resolution.

555

556 Several recent studies suggest that the IMERG represents a substantial advancement  
557 over TRMM and likely the others (e.g., Wei et al. 2017; Khodadoust Siuki et al. 2017;  
558 Zhang et al. 2018), thus we rely on IMERG as the default in much of our analysis.  
559 However, we do not entirely discount the other products because the discrepancy  
560 between them provides a measure of uncertainty in the satellite-based estimates of  
561 precipitation. Our use of the minimum to maximum range of multiple observational  
562 products is indicative of their discrepancy, but not their uncertainty, and thus is a limitation  
563 of the current work and challenge that we hope will be addressed in the future.

564

565 The common model biases identified in this study are mainly associated with the  
566 overestimated light precipitation and underestimated heavy precipitation. These biases  
567 persist from deficiencies identified in earlier generation models (e.g., Dai 2006), and as  
568 shown in this study there has been little improvement. One reason may be that these key  
569 characteristics of precipitation are not commonly considered in the model development

570 process. Enabling modelers to more readily objectively evaluate simulated precipitation  
571 distributions could perhaps serve as a guide to improvement. The current study aims to  
572 provide a framework for objective evaluation of simulated precipitation distributions at  
573 regional scales.

574

575 Imperfect convective parameterizations are a possible cause of the common model  
576 biases in precipitation distributions (e.g., Lin et al. 2013; Kooperman et al. 2018; Ahn et  
577 al. 2018; Chen and Dai 2019; Chen et al. 2021; Martinez-Villalobos et al. 2022). Many  
578 convective parameterizations tend to produce too frequent and light precipitation, the so-  
579 called “drizzling” bias (e.g., Dai 2006; Trenberth et al. 2017; Chen et al. 2021; Ma et al.  
580 2022), and it is likely due to a fact that the parameterized convection is more readily  
581 triggered than that in the nature (e.g., Lin et al. 2013; Chen et al. 2021). As model  
582 horizontal resolution increases, grid-scale precipitation processes can lead to resolving  
583 convective precipitation, as in so-called cloud resolving, storm resolving, or convective  
584 permitting models. Ma et al. (2022) compare several storm resolving models in  
585 DYAMOND to recent CMIP6 models with a convective parameterization and observe that  
586 the simulated precipitation distributions are more realistic in the storm resolving models.  
587 However, some of the storm resolving models still suffer from precipitation distribution  
588 errors, including bimodality in the frequency distribution. Further studies are needed to  
589 better understand the precipitation distribution biases in models.

590

591

592 **6. Conclusion**

593 We introduce a framework for regional scale evaluation of simulated precipitation  
594 distributions with 62 climate reference regions and 10 precipitation distribution metrics  
595 and apply it to evaluate the two most recent generations of climate model intercomparison  
596 simulations (i.e., CMIP5 and CMIP6).

597

598 To facilitate the regional scale for evaluation, regions where precipitation characteristics  
599 are relatively homogenous are identified. Our reference regions consist of existing IPCC  
600 AR6 climate reference regions, with additional subdivisions based on homogeneity  
601 analysis performed on precipitation distributions within each region. Our precipitation  
602 clustering analysis reveals that the IPCC AR6 land regions are reasonably homogeneous  
603 in precipitation character, while some ocean regions are relatively inhomogeneous,  
604 including large portions of both heavy and light precipitating areas. To define more  
605 homogeneous regions for the analysis of precipitation distributions, we have modified  
606 some ocean regions to better fit the clustering results. Although the clustering regions are  
607 obtained based on the IMERG annual precipitation, the improved homogeneity is fairly  
608 consistent across different datasets (TRMM, CMORPH, GPCP, PERSIANN, and ERA5)  
609 and seasons (MAM, JJA, SON, and DJF). Use of these more homogeneous regions  
610 enables us to extract more robust quantitative information from the distributions in each  
611 region.

612

613 To form the basis for evaluation within each region, we use a set of metrics that are well-  
614 established and easy to interpret, aiming to extract key characteristics from the  
615 distributions of precipitation frequency, amount, and cumulative fraction of precipitation

616 amount. We include the precipitation rate at the peak of the amount and frequency  
617 distributions (Kooperman et al., 2016; Pendergrass and Deser, 2017) and define several  
618 complementary metrics to measure the frequency and amount of precipitation under the  
619 10th percentile (P10) and over the 90th percentile (P90). The distribution peak metrics  
620 assess whether the center of each distribution is shifted toward light or heavy  
621 precipitation, while the P10 and P90 metrics quantify the fraction of light and heavy  
622 precipitation in the distributions. The Perkins score is included to measure the similarity  
623 between the observed and modeled frequency distributions. Also, based on the  
624 cumulative fraction of precipitation amount, we implement the unevenness metric  
625 counting the number of wettest days for half of the annual precipitation (Pendergrass and  
626 Knutti 2018), the fraction of annual precipitating days above 1 mm/day, and the simple  
627 daily intensity index (Zhang et al. 2011).

628

629 We apply the framework of regional scale precipitation distribution benchmarking to all  
630 available realizations of 25 CMIP5 and 41 CMIP6 models and 5 satellite-based  
631 precipitation products (IMERG, TRMM, CMORPH, GPCP, PERSIANN). The  
632 observational discrepancy is substantially larger compared to the models' spread for  
633 some regions, especially for mid-latitude and polar regions and for some metrics such as  
634 Amount P90 and Frequency P90. We use two approaches to account for observational  
635 discrepancy in the model evaluation. One is based on the number of models within the  
636 observational range, and another is the number of models below/above all observations.  
637 In this way, we can draw some conclusions on the overall performance in the CMIP  
638 ensemble even in the presence of observations that may substantially disagree in certain

639 regions. Many CMIP5 and CMIP6 models underestimate the Amount and Frequency  
640 peaks and overestimate Amount and Frequency P10 compared to observations,  
641 especially in many mid-latitude regions where more than 50% of the models are out of  
642 the observational range. This indicates that models produce too frequent light  
643 precipitation, a bias that is also revealed by the overestimated FracPRdays and the  
644 underestimated SDII. Unevenness is the metric that models simulate the worst – in many  
645 regions more than 70-90% of the models are out of the observational range. Clear  
646 changes in performance between CMIP5 and CMIP6 are limited. Considering all metrics,  
647 the CMIP6 models show improvement in some mid-latitude regions, but in a few tropical  
648 regions the CMIP6 models actually show performance degradation.

649

650 The framework presented in this study is intended to be a useful resource for model  
651 evaluation analysts and developers working towards improved performance for a wide  
652 range of precipitation characteristics. Basing the regions in part on homogeneous  
653 precipitation characteristics can facilitate identification of the processes responsible for  
654 model errors as heavy precipitating regions are generally dominated by convective  
655 precipitation, while the moderate and light precipitation regions are mainly governed by  
656 stratiform precipitation processes. Although the framework presented herein has been  
657 demonstrated with regional scale evaluation benchmarking, it can be applicable for  
658 benchmarking at larger scales and homogeneous precipitation regions.

659

660 **Code Availability**

661 The benchmarking framework for precipitation distributions established in this study is  
662 available via the PCMDI Metrics Package (PMP,  
663 [https://github.com/PCMDI/pcmdi\\_metrics](https://github.com/PCMDI/pcmdi_metrics), DOI: [10.5281/zenodo.7231033](https://doi.org/10.5281/zenodo.7231033)). This  
664 framework provides three tiers of area averaged outputs for i) large scale domain (Tropics  
665 and Extratropics with separated land and ocean) commonly used in the PMP, ii) large  
666 scale domain with clustered precipitation characteristics (Tropics and Extratropics with  
667 separated land and ocean, and separated heavy, moderate, and light precipitation  
668 regions), and iii) modified IPCC AR6 regions shown in this paper.

669

670

671 **Data Availability**

672 All of the data used in this study are publicly available. The satellite-based precipitation  
673 products used in this study (IMERG, TRMM, CMORPH, GPCP, and PERSIANN) and  
674 ERA5 precipitation product are available on the Obs4MIPs at [https://esgf-](https://esgf-node.llnl.gov/projects/obs4mips/)  
675 [node.llnl.gov/projects/obs4mips/](https://esgf-node.llnl.gov/projects/obs4mips/). The CMIP data is available on the ESGF at [https://esgf-](https://esgf-node.llnl.gov/projects/esgf-llnl)  
676 [node.llnl.gov/projects/esgf-llnl](https://esgf-node.llnl.gov/projects/esgf-llnl). The statistics generated from this benchmarking  
677 framework and the interactive plots with access to the underlying diagnostics were made  
678 available on the PCMDI Simulation Summaries at  
679 <https://pcmdi.llnl.gov/research/metrics/precip/>.

680

681

682 **Author contribution**

683 PG and AP designed the initial idea of the precipitation benchmarking framework. MA,  
684 PU, PG, and JL advanced the idea and developed the framework. MA performed  
685 analysis. MA, JL, and AO implemented the framework code into the PCMDI metrics  
686 package. MA prepared the manuscript with contributions from all co-authors.

687

688

### 689 **Competing interests**

690 The authors declare that they have no conflict of interest.

691

692

### 693 **Disclaimer**

694 This document was prepared as an account of work sponsored by an agency of the U.S.  
695 government. Neither the U.S. government nor Lawrence Livermore National Security,  
696 LLC, nor any of their employees makes any warranty, expressed or implied, or assumes  
697 any legal liability or responsibility for the accuracy, completeness, or usefulness of any  
698 information, apparatus, product, or process disclosed, or represents that its use would  
699 not infringe privately owned rights. Reference herein to any specific commercial product,  
700 process, or service by trade name, trademark, manufacturer, or otherwise does not  
701 necessarily constitute or imply its endorsement, recommendation, or favoring by the U.S.  
702 government or Lawrence Livermore National Security, LLC. The views and opinions of  
703 authors expressed herein do not necessarily state or reflect those of the U.S. government  
704 or Lawrence Livermore National Security, LLC, and shall not be used for advertising or  
705 product endorsement purposes.



706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723

**Acknowledgements**

This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344. The efforts of the authors were supported by the Regional and Global Model Analysis (RGMA) program of the United States Department of Energy's Office of Science, including under Award Number DE-SC0022070 and National Science Foundation (NSF) IA 1947282. This work was also partially supported by the National Center for Atmospheric Research (NCAR), which is a major facility sponsored by the NSF under Cooperative Agreement No. 1852977. We acknowledge the World Climate Research Programme's Working Group on Coupled Modeling, which is responsible for CMIP, and we thank the climate modeling groups for producing and making available their model output, the Earth System Grid Federation (ESGF) for archiving the output and providing access, and the multiple funding agencies who support CMIP and ESGF. The U.S. Department of Energy's Program for Climate Model Diagnosis and Intercomparison (PCMDI) provides coordinating support and led development of software infrastructure for CMIP.

724 **References**

725 Abramowitz, G. (2012). Towards a public, standardized, diagnostic benchmarking  
726 system for land surface models. *Geoscientific Model Development*, 5(3), 819–  
727 827. <https://doi.org/10.5194/gmd-5-819-2012>.

728 Ahn, M., and I. Kang, 2018: A practical approach to scale-adaptive deep convection  
729 in a GCM by controlling the cumulus base mass flux. *npj Clim. Atmos. Sci.*, 1,  
730 13, <https://doi.org/10.1038/s41612-018-0021-0>.

731 Ahn, M.-S., P. A. Ullrich, J. Lee, P. J. Gleckler, H.-Y. Ma, C. R. Terai, P. A.  
732 Bogenschutz, and A. C. Ordonez, 2023: Bimodality in Simulated Precipitation  
733 Frequency Distributions and Its Relationship with Convective Parameterizations.  
734 *npj Climate and Atmospheric Science*, submitted.

735 Ahn, M.-S., P. J. Gleckler, J. Lee, A. G. Pendergrass, and C. Jakob, 2022:  
736 Benchmarking Simulated Precipitation Variability Amplitude across Time  
737 Scales. *J. Clim.*, **35**, 3173–3196, <https://doi.org/10.1175/JCLI-D-21-0542.1>.

738 Ashouri, H., K. L. Hsu, S. Sorooshian, D. K. Braithwaite, K. R. Knapp, L. D. Cecil, B.  
739 R. Nelson, and O. P. Prat, 2015: PERSIANN-CDR: Daily precipitation climate  
740 data record from multisatellite observations for hydrological and climate studies.  
741 *Bull. Am. Meteorol. Soc.*, **96**, 69–83, [https://doi.org/10.1175/BAMS-D-13-](https://doi.org/10.1175/BAMS-D-13-00068.1)  
742 00068.1.

743 Chakravarti, I. M., R. G. Laha, and J. Roy, 1967: Handbook of Methods of Applied  
744 Statistics, Volume I: Techniques of Computation, Descriptive Methods, and  
745 Statistical Inference. *John Wiley Sons*, 392–394.

746 Chen, D., and A. Dai, 2019: Precipitation Characteristics in the Community  
747 Atmosphere Model and Their Dependence on Model Physics and Resolution. *J.*  
748 *Adv. Model. Earth Syst.*, **11**, 2352–2374,  
749 <https://doi.org/10.1029/2018MS001536>.

750 Chen, D., A. Dai, and A. Hall, 2021: The Convective-To-Total Precipitation Ratio and  
751 the “Drizzling” Bias in Climate Models. *J. Geophys. Res. Atmos.*, **126**, 1–17,  
752 <https://doi.org/10.1029/2020JD034198>.

753 Covey, C., Gleckler, P. J., Doutriaux, C., Williams, D. N., Dai, A., Fasullo, J.,  
754 Trenberth, K., & Berg, A. (2016). Metrics for the Diurnal Cycle of Precipitation:  
755 Toward Routine Benchmarks for Climate Models. *Journal of Climate*, *29*(12),  
756 4461–4471. <https://doi.org/10.1175/JCLI-D-15-0664.1>

757 Dai, A., 2006: Precipitation characteristics in eighteen coupled climate models. *J.*  
758 *Clim.*, **19**, 4605–4630, <https://doi.org/10.1175/JCLI3884.1>.

759 Eyring, V., S. Bony, G. A. Meehl, C. A. Senior, B. Stevens, R. J. Stouffer, and K. E.  
760 Taylor, 2016: Overview of the Coupled Model Intercomparison Project Phase 6  
761 (CMIP6) experimental design and organization. *Geosci. Model Dev.*, **9**, 1937–  
762 1958, <https://doi.org/10.5194/gmd-9-1937-2016>.

763 Fiedler, S., and Coauthors, 2020: Simulated Tropical Precipitation Assessed across  
764 Three Major Phases of the Coupled Model Intercomparison Project (CMIP).  
765 *Mon. Weather Rev.*, **148**, 3653–3680, [https://doi.org/10.1175/MWR-D-19-](https://doi.org/10.1175/MWR-D-19-0404.1)  
766 0404.1.

767 Gleckler, P., C. Doutriaux, P. Durack, K. Taylor, Y. Zhang, D. Williams, E. Mason,  
768 and J. Servonnat, 2016: A More Powerful Reality Test for Climate Models. *Eos*  
769 (*Washington, DC*), **97**, 20–24, <https://doi.org/10.1029/2016EO051663>.

770 Gleckler, P. J., K. E. Taylor, and C. Doutriaux, 2008: Performance metrics for climate  
771 models. *J. Geophys. Res. Atmos.*, **113**, 1–20,  
772 <https://doi.org/10.1029/2007JD008972>.

773 Hersbach, H., and Coauthors, 2020: The ERA5 global reanalysis. *Q. J. R. Meteorol.*  
774 *Soc.*, **146**, 1999–2049, <https://doi.org/10.1002/qj.3803>.

775 Huffman, G. J., and Coauthors, 2007: The TRMM Multisatellite Precipitation Analysis  
776 (TMPA): Quasi-Global, Multiyear, Combined-Sensor Precipitation Estimates at  
777 Fine Scales. *J. Hydrometeorol.*, **8**, 38–55, <https://doi.org/10.1175/JHM560.1>.

778 Huffman, G. J., and Coauthors, 2020: Integrated Multi-satellite Retrievals for the  
779 Global Precipitation Measurement (GPM) Mission (IMERG). *Advances in Global*  
780 *Change Research*, Vol. 67 of, 343–353.

781 Huffman, G. J., R. F. Adler, M. M. Morrissey, D. T. Bolvin, S. Curtis, R. Joyce, B.  
782 McGavock, and J. Susskind, 2001: Global Precipitation at One-Degree Daily

783 Resolution from Multisatellite Observations. *J. Hydrometeorol.*, **2**, 36–50,  
784 [https://doi.org/10.1175/1525-7541\(2001\)002<0036:GPAODD>2.0.CO;2](https://doi.org/10.1175/1525-7541(2001)002<0036:GPAODD>2.0.CO;2).

785 Iturbide, M., and Coauthors, 2020: An update of IPCC climate reference regions for  
786 subcontinental analysis of climate model data: definition and aggregated  
787 datasets. *Earth Syst. Sci. Data*, **12**, 2959–2970, [https://doi.org/10.5194/essd-12-](https://doi.org/10.5194/essd-12-2959-2020)  
788 [2959-2020](https://doi.org/10.5194/essd-12-2959-2020).

789 Khodadoust Siuki, S., B. Saghafian, and S. Moazami, 2017: Comprehensive  
790 evaluation of 3-hourly TRMM and half-hourly GPM-IMERG satellite precipitation  
791 products. *Int. J. Remote Sens.*, **38**, 558–571,  
792 <https://doi.org/10.1080/01431161.2016.1268735>.

793 Kim, S., A. Sharma, C. Wasko, and R. Nathan, 2022: Linking Total Precipitable Water  
794 to Precipitation Extremes Globally. *Earth's Futur.*, **10**,  
795 <https://doi.org/10.1029/2021EF002473>.

796 Kooperman, G. J., M. S. Pritchard, M. A. Burt, M. D. Branson, and D. A. Randall,  
797 2016: Robust effects of cloud superparameterization on simulated daily rainfall  
798 intensity statistics across multiple versions of the Community Earth System  
799 Model. *J. Adv. Model. Earth Syst.*, **8**, 140–165,  
800 <https://doi.org/10.1002/2015MS000574>.

801 Kooperman, G. J., M. S. Pritchard, T. A. O'Brien, and B. W. Timmermans, 2018:  
802 Rainfall From Resolved Rather Than Parameterized Processes Better  
803 Represents the Present-Day and Climate Change Response of Moderate Rates

804 in the Community Atmosphere Model. *J. Adv. Model. Earth Syst.*, **10**, 971–988,  
805 <https://doi.org/10.1002/2017MS001188>.

806 Leung, L. R., and Coauthors, 2022: Exploratory Precipitation Metrics: Spatiotemporal  
807 Characteristics, Process-Oriented, and Phenomena-Based. *J. Clim.*, **35**, 3659–  
808 3686, <https://doi.org/10.1175/JCLI-D-21-0590.1>.

809 Lin, Y., M. Zhao, Y. Ming, J.-C. Golaz, L. J. Donner, S. A. Klein, V. Ramaswamy, and  
810 S. Xie, 2013: Precipitation Partitioning, Tropical Clouds, and Intraseasonal  
811 Variability in GFDL AM2. *J. Clim.*, **26**, 5453–5466, <https://doi.org/10.1175/JCLI-D-12-00442.1>.

813 Ma, H., S. A. Klein, J. Lee, M. Ahn, C. Tao, and P. J. Gleckler, 2022: Superior Daily  
814 and Sub-Daily Precipitation Statistics for Intense and Long-Lived Storms in  
815 Global Storm-Resolving Models. *Geophys. Res. Lett.*, **49**,  
816 <https://doi.org/10.1029/2021GL096759>.

817 MacQueen, J. B., 1967: Some methods for classification and analysis of multivariate  
818 observations. *Berkeley Symp. Math. Stat. Probab.*, **VOL. 5.1**, 281–297.

819 Martinez-Villalobos, C., and J. D. Neelin, 2019: Why Do Precipitation Intensities Tend  
820 to Follow Gamma Distributions? *J. Atmos. Sci.*, **76**, 3611–3631,  
821 <https://doi.org/10.1175/JAS-D-18-0343.1>.

822 Martinez-Villalobos, C., J. D. Neelin, and A. G. Pendergrass, 2022: Metrics for  
823 Evaluating CMIP6 Representation of Daily Precipitation Probability Distributions.  
824 *J. Clim.*, 1–79, <https://doi.org/10.1175/JCLI-D-21-0617.1>.

825 Meehl, G. A., C. Covey, B. McAvaney, M. Latif, and R. J. Stouffer, 2005: Overview  
826 of the Coupled Model Intercomparison Project. *Bull. Am. Meteorol. Soc.*, **86**, 89–  
827 96, <https://doi.org/10.1175/BAMS-86-1-89>.

828 Meehl, G. A., C. Covey, T. Delworth, M. Latif, B. McAvaney, J. F. B. Mitchell, R. J.  
829 Stouffer, and K. E. Taylor, 2007: THE WCRP CMIP3 Multimodel Dataset: A New  
830 Era in Climate Change Research. *Bull. Am. Meteorol. Soc.*, **88**, 1383–1394,  
831 <https://doi.org/10.1175/BAMS-88-9-1383>.

832 Meehl, G. A., G. J. Boer, C. Covey, M. Latif, and R. J. Stouffer, 2000: The Coupled  
833 Model Intercomparison Project (CMIP). *Bull. Am. Meteorol. Soc.*, **81**, 313–318,  
834 [https://doi.org/10.1175/1520-0477\(2000\)081<0313:TCMIPC>2.3.CO;2](https://doi.org/10.1175/1520-0477(2000)081<0313:TCMIPC>2.3.CO;2).

835 Pendergrass, A. G., and C. Deser, 2017: Climatological Characteristics of Typical  
836 Daily Precipitation. *J. Clim.*, **30**, 5985–6003, [https://doi.org/10.1175/JCLI-D-16-](https://doi.org/10.1175/JCLI-D-16-0684.1)  
837 [0684.1](https://doi.org/10.1175/JCLI-D-16-0684.1).

838 Pendergrass, A. G., and D. L. Hartmann, 2014: Two Modes of Change of the  
839 Distribution of Rain\*. *J. Clim.*, **27**, 8357–8371, [https://doi.org/10.1175/JCLI-D-](https://doi.org/10.1175/JCLI-D-14-00182.1)  
840 [14-00182.1](https://doi.org/10.1175/JCLI-D-14-00182.1).

841 Pendergrass, A. G., and R. Knutti, 2018: The Uneven Nature of Daily Precipitation  
842 and Its Change. *Geophys. Res. Lett.*, **45**, 11,980–11,988,  
843 <https://doi.org/10.1029/2018GL080298>.

844 Pendergrass, A. G., P. J. Gleckler, L. R. Leung, and C. Jakob, 2020: Benchmarking  
845 Simulated Precipitation in Earth System Models. *Bull. Am. Meteorol. Soc.*, **101**,  
846 E814–E816, <https://doi.org/10.1175/BAMS-D-19-0318.1>.

847 Perkins, S. E., A. J. Pitman, N. J. Holbrook, and J. McAneney, 2007: Evaluation of  
848 the AR4 Climate Models' Simulated Daily Maximum Temperature, Minimum  
849 Temperature, and Precipitation over Australia Using Probability Density  
850 Functions. *J. Clim.*, **20**, 4356–4376, <https://doi.org/10.1175/JCLI4253.1>.

851 Roca, R., L. V. Alexander, G. Potter, M. Bador, R. Jucá, S. Contractor, M. G.  
852 Bosilovich, and S. Cloché, 2019: FROGS: a daily 1° × 1° gridded precipitation  
853 database of rain gauge, satellite and reanalysis products. *Earth Syst. Sci. Data*,  
854 **11**, 1017–1035, <https://doi.org/10.5194/essd-11-1017-2019>.

855 Stephens, M. A., 1974: EDF Statistics for Goodness of Fit and Some Comparisons.  
856 *J. Am. Stat. Assoc.*, **69**, 730–737, <https://doi.org/10.2307/2286009>.

857 Sun, Y., S. Solomon, A. Dai, and R. W. Portmann, 2006: How Often Does It Rain? *J.*  
858 *Clim.*, **19**, 916–934, <https://doi.org/10.1175/JCLI3672.1>.

859 Sun, Y., S. Solomon, A. Dai, and R. W. Portmann, 2007: How Often Will It Rain? *J.*  
860 *Clim.*, **20**, 4801–4818, <https://doi.org/10.1175/JCLI4263.1>.

861 Swenson, L. M., and R. Grotjahn, 2019: Using Self-Organizing Maps to Identify  
862 Coherent CONUS Precipitation Regions. *J. Clim.*, **32**, 7747–7761,  
863 <https://doi.org/10.1175/JCLI-D-19-0352.1>.



864 Tang, S., P. Gleckler, S. Xie, J. Lee, M.-S. Ahn, C. Covey, and C. Zhang, 2021:  
865 Evaluating Diurnal and Semi-Diurnal Cycle of Precipitation in CMIP6 Models  
866 Using Satellite- and Ground-Based Observations. *J. Clim.*, 1–56,  
867 <https://doi.org/10.1175/JCLI-D-20-0639.1>.

868 Taylor, K. E., R. J. Stouffer, and G. A. Meehl, 2012: An overview of CMIP5 and the  
869 experiment design. *Bull. Amer. Meteor. Soc.*, **93**, 485–498,  
870 <https://doi.org/10.1175/BAMS-D-11-00094.1>.

871 Trenberth, K. E., A. Dai, R. M. Rasmussen, and D. B. Parsons, 2003: The Changing  
872 Character of Precipitation. *Bull. Am. Meteorol. Soc.*, **84**, 1205–1218,  
873 <https://doi.org/10.1175/BAMS-84-9-1205>.

874 Trenberth, K. E., and Y. Zhang, 2018: How Often Does It Really Rain? *Bull. Am.*  
875 *Meteorol. Soc.*, **99**, 289–298, <https://doi.org/10.1175/BAMS-D-17-0107.1>.

876 Trenberth, K. E., Y. Zhang, and M. Gehne, 2017: Intermittency in Precipitation:  
877 Duration, Frequency, Intensity, and Amounts Using Hourly Data. *J.*  
878 *Hydrometeorol.*, **18**, 1393–1412, <https://doi.org/10.1175/JHM-D-16-0263.1>.

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

883 Waliser, D., and Coauthors, 2020: Observations for Model Intercomparison Project  
884 (Obs4MIPs): status for CMIP6. *Geosci. Model Dev.*, **13**, 2945–2958,  
885 <https://doi.org/10.5194/gmd-13-2945-2020>.

886 Wehner, M., P. Gleckler, J. Lee, 2020: Characterization of long period return values  
887 of extreme daily temperature and precipitation in the CMIP6 models: Part 1,  
888 model evaluation. *Weather and Climate Extremes*, **30**, 100283, doi:  
889 [10.1016/j.wace.2020.100283](https://doi.org/10.1016/j.wace.2020.100283).

890 Wei, G., H. Lü, W. T. Crow, Y. Zhu, J. Wang, and J. Su, 2017: Evaluation of Satellite-  
891 Based Precipitation Products from IMERG V04A and V03D, CMORPH and  
892 TMPA with Gauged Rainfall in Three Climatologic Zones in China. *Remote*  
893 *Sens.*, **10**, 30, <https://doi.org/10.3390/rs10010030>.

894 Xie, P., R. Joyce, S. Wu, S. H. Yoo, Y. Yarosh, F. Sun, and R. Lin, 2017:  
895 Reprocessed, bias-corrected CMORPH global high-resolution precipitation  
896 estimates from 1998. *J. Hydrometeorol.*, **18**, 1617–1641,  
897 <https://doi.org/10.1175/JHM-D-16-0168.1>.

898 Zhang, C., X. Chen, H. Shao, S. Chen, T. Liu, C. Chen, Q. Ding, and H. Du, 2018:  
899 Evaluation and intercomparison of high-resolution satellite precipitation  
900 estimates-GPM, TRMM, and CMORPH in the Tianshan Mountain Area. *Remote*  
901 *Sens.*, **10**, <https://doi.org/10.3390/rs10101543>.

902 Zhang, X., L. Alexander, G. C. Hegerl, P. Jones, A. K. Tank, T. C. Peterson, B.  
903 Trewin, and F. W. Zwiers, 2011: Indices for monitoring changes in extremes

904 based on daily temperature and precipitation data. *Wiley Interdiscip. Rev. Clim.*  
905 *Chang.*, **2**, 851–870, <https://doi.org/10.1002/wcc.147>.

906 **Tables**

907

908

909

910 Table 1. Satellite-based and reanalysis precipitation products used in this study.

911

Product	Data source	Coverage		Resolution		Reference
		Domain	Period	Horizontal	Frequency	
IMERG	NASA Integrated Multi-satellite Retrievals for GPM version 6 final run product	Global, while beyond 60°NS is incomplete	2000.6-present	0.1°	30 minutes	Huffman et al. (2020)
TRMM	NASA Tropical Rainfall Measuring Mission Multi-satellite Precipitation Analysis 3B42 version 7 product	50°S-50°N	1998.1-2019.12	0.25°	3 hours	Huffman et al. (2007)
CMORPH	NOAA Bias-corrected Climate Prediction Center Morphing technique product	60°S-60°N	1998.1-present	0.073°	30 minutes	Xie et al. (2017)
GPCP	NASA Global Precipitation Climatology Project 1DD version 1.3	Global, while beyond 40°NS is incomplete	1996.10-present	1°	1 day	Huffman et al. (2001)
PERSIANN	UC-IRVINE/CHRS Precipitation Estimation from Remotely Sensed Information using Artificial Neural Networks-Climate Data Record	60°S-60°N	1983.1-present	0.25°	1 day	Ashouri et al. (2015)
ERA5	ECMWF Integrated Forecasting System Cy41r2	Global	1950.1-present	0.25°	1 hour	Hersbach et al. (2020)

912

913

914

915

916

917

918

919

920

921

922

923

924

925  
926  
927  
928  
929  
930

Table 2. CMIP5 and CMIP6 models used in this study and their horizontal resolution. The number in parentheses indicates the number of realizations used for each model. Note that the horizontal resolution information is obtained from the number of grids, and it may vary slightly if the grid interval is not linear.

Institute	CMIP5		CMIP6	
	Name	Horizontal resolution [lon x lat °]	Name	Horizontal resolution [lon x lat °]
CSIRO/BOM, Australia	ACCESS1-0 (1)	1.875 x 1.241	ACCESS-CM2 (7)	1.875 x 1.25
	ACCESS1-3 (2)	1.875 x 1.241	ACCESS-ESM1-5 (10)	1.875 x 1.241
BCC, China	BCC-CSM1-1 (3)	1.875 x 1.241	BCC-CSM2-MR (3)	1.125 x 1.125
	BCC-CSM1-1-M (3)	1.125 x 1.125	BCC-ESM1 (3)	2.812 x 2.812
BNU, China	BNU-ESM (1)	2.812 x 2.812	N/A	
CAMS, China	N/A		CAMS-CSM1-0 (3)	
CCCma, Canada	N/A		CanESM5 (7)	2.812 x 2.812
NCAR, USA	CCSM4 (6)	1.25 x 0.938	CESM2 (10)	1.25 x 0.938
			CESM2-FV2 (3)	2.5 x 1.875
			CESM2-WACCM (3)	1.25 x 0.938
			CESM2-WACCM-FV2 (3)	2.5 x 1.875
CMCC, Italy	CMCC-CM (3)	0.75 x 0.75	CMCC-CM2-HR4 (1)	1.25 x 0.938
			CMCC-CM2-SR5 (1)	1.25 x 0.938
CNRM-CERFACS, France	N/A		CNRM-CM6-1 (1)	1.406 x 1.406
			CNRM-CM6-1-HR (1)	0.5 x 0.5
			CNRM-ESM2-1 (1)	1.406 x 1.406
CSIRO-QCCCE, Australia	CSIRO-Mk3-6-0 (10)	1.875 x 1.875	N/A	
DOE, USA	N/A		E3SM-1-0 (3)	1.0 x 1.0
EC-Earth Consortium, European Community	EC-Earth (1)	1.125 x 1.125	EC-Earth3 (6)	0.703 x 0.703
			EC-Earth3-AerChem (1)	0.703 x 0.703
			EC-Earth3-CC (5)	
			EC-Earth3-Veg (3)	0.703 x 0.703
IAP-CAS/THU, China	FGOALS-g2 (1)	2.812 x 3.0	FGOALS-f3-L (3)	1.0 x 1.0
	FGOALS-s2 (3)	2.812 x 1.667		
NOAA GFDL, USA	GFDL-CM3 (5)	2.5 x 2.0	GFDL-CM4 (1)	1.0 x 1.0
	GFDL-HIRAM-C180 (2)	0.625 x 0.5	GFDL-ESM4 (1)	1.0 x 1.0
	GFDL-HIRAM-C360 (1)	0.312 x 0.25		
NASA GISS, USA	GISS-E2-R (2)	2.5 x 2.0	N/A	
MOHC, UK	HadGEM2-A (1)	1.875 x 1.241	HadGEM3-GC31-LL (5)	1.875 x 1.25
			HadGEM3-GC31-MM (4)	0.833 x 0.556
			UKESM1-0-LL (1)	1.875 x 1.25
IITM, India	N/A		IITM-ESM (1)	1.875 x 1.915
INM, Russia	INMCM4 (1)	2.0 x 1.5	INM-CM4-8 (1)	2.0 x 1.5
			INM-CM5-0 (1)	2.0 x 1.5

IPSL, France	IPSL-CM5A-LR (6)	3.75 x 1.875	IPSL-CM6A-LR (22)	2.5 x 1.259
	IPSL-CM5A-MR (3)	2.5 x 1.259		
	IPSL-CM5B-LR (1)	3.75 x 1.875		
NIMS/KMA, Korea	N/A		KACE-1-0-G (1)	1.875 x 1.25
MIROC, Japan	MIROC5 (2)	1.406 x 1.406	MIROC6 (10)	1.406 x 1.406
			MIROC-ES2L (3)	2.812 x 2.812
MPI-M, Germany	MPI-ESM-MR (3)	1.875 x 1.875	MPI-ESM-1-2-HAM (3)	1.875 x 1.875
			MPI-ESM1-2-HR (3)	0.938 x 0.938
			MPI-ESM1-2-LR (3)	1.875 x 1.875
MRI, Japan	MRI-AGCM3-2H (1)	0.562 x 0.562	MRI-ESM2-0 (3)	1.125 x 1.125
	MRI-AGCM3-2S (1)	0.188 x 0.188		
	MRI-CGCM3 (3)	1.125 x 1.125		
NCC, Norway	N/A		NorCPM1 (10)	2.5 x 1.875
			NorESM2-LM (2)	2.5 x 1.875
SNU, Korea	N/A		SAM0-UNICON (1)	1.25 x 0.938
AS-RCEC, Taiwan	N/A		TaiESM1 (1)	1.25 x 0.938

931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955

956  
957  
958

Table 3. Precipitation distribution metrics implemented in this study.

<b>Metric [unit]</b>	<b>Definition</b>	<b>Objectives</b>	<b>Reference</b>
<b>Amount peak</b> [mm/day]	Rain rate where the maximum rain amount occurs	Characterize typical daily precipitation amount	Pendergrass and Deser (2017)
<b>Amount P10</b> [fraction]	Fraction of rain amount in lower 10 percentile of OBS amount	Measure the rain amount from light rainfall	
<b>Amount P90</b> [fraction]	Fraction of rain amount in upper 90 percentile of OBS amount	Measure the rain amount from heavy rainfall	
<b>Frequency peak</b> [mm/day]	Rain rate where the maximum nonzero rain frequency occurs	Characterize typical daily precipitation frequency	Pendergrass and Deser (2017)
<b>Frequency P10</b> [fraction]	Fraction of rain frequency in lower 10 percentile of OBS amount	Measure the frequency of light rainfall	
<b>Frequency P90</b> [fraction]	Fraction of rain frequency in upper 90 percentile of OBS amount	Measure the frequency of heavy rainfall	
<b>Unevenness</b> [days]	Number of wettest days for that constitute half of annual precipitation	Measure uneven characteristic of daily precipitation	Pendergrass and Knutti (2018)
<b>FracPRdays</b> [fraction]	Number of precipitating days ( $\geq 1$ mm/day) divided by total days a year	Measure fraction of precipitating days a year	Updated from Zhang et al. (2011)
<b>SDII</b> [mm/day]	Annual total precipitation divided by the number of precipitating days ( $\geq 1$ mm/day)	Measure daily precipitation intensity	Zhang et al. (2011)
<b>Perkins score</b> [unitless between 0-1]	Sum of minimum values between two PDFs across all bins	Measure similarity between two PDFs	Perkins et al. (2007)

959  
960  
961  
962  
963  
964

965  
966  
967  
968

Table 4. List of climate reference regions used in this study. The new ocean regions defined in this study are highlighted in bold.

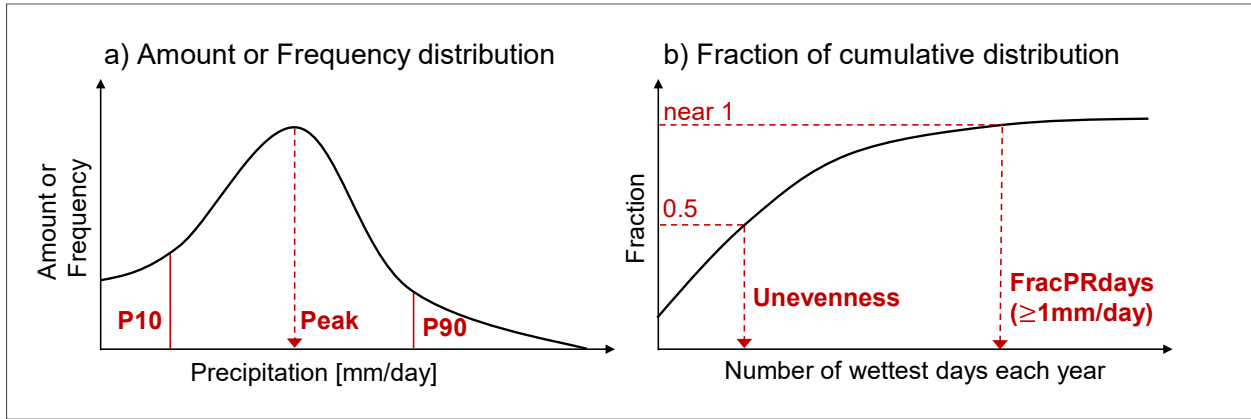
1	GIC	Greenland/Iceland	22	WAF	Western-Africa	43	SAU	S.Australia
2	NWN	N.W.North-America	23	CAF	Central-Africa	44	NZ	New-Zealand
3	NEN	N.E.North-America	24	NEAF	N.Eastern-Africa	45	EAN	E.Antarctica
4	WNA	W.North-America	25	SEAF	S.Eastern-Africa	46	WAN	W.Antarctica
5	CNA	C.North-America	26	WSAF	W.Southern-Africa	47	ARO	Arctic-Ocean
6	ENA	E.North-America	27	ESAF	E.Southern-Africa	48	ARS	Arabian-Sea
7	NCA	N.Central-America	28	MDG	Madagascar	49	BOB	Bay-of-Bengal
8	SCA	S.Central-America	29	RAR	Russian-Arctic	50	EIO	Equatorial-Indian-Ocean
9	CAR	Caribbean	30	WSB	W.Siberia	51	SIO	S.Indian-Ocean
10	NWS	N.W.South-America	31	ESB	E.Siberia	52	<b>NPO</b>	<b>N.Pacific-Ocean</b>
11	NSA	N.South-America	32	RFE	Russian-Far-East	53	<b>NWP O</b>	<b>N.W.Pacific-Ocean</b>
12	NES	N.E.South-America	33	WCA	W.C.Asia	54	<b>NEPO</b>	<b>N.E.Pacific-Ocean</b>
13	SAM	South-American-Monsoon	34	ECA	E.C.Asia	55	<b>PITCZ</b>	<b>Pacific-ITCZ</b>
14	SWS	S.W.South-America	35	TIB	Tibetan-Plateau	56	<b>SWPO</b>	<b>S.W.Pacific-Ocean</b>
15	SES	S.E.South-America	36	EAS	E.Asia	57	<b>SEPO</b>	<b>S.E.Pacific-Ocean</b>
16	SSA	S.South-America	37	ARP	Arabian-Peninsula	58	<b>NAO</b>	<b>N.Atlantic-Ocean</b>
17	NEU	N.Europe	38	SAS	S.Asia	59	<b>NEAO</b>	<b>N.E.Atlantic-Ocean</b>
18	WCE	West&Central-Europe	39	SEA	S.E.Asia	60	<b>AITCZ</b>	<b>Atlantic-ITCZ</b>
19	EEU	E.Europe	40	NAU	N.Australia	61	<b>SAO</b>	<b>S.Atlantic-Ocean</b>
20	MED	Mediterranean	41	CAU	C.Australia	62	<b>SOO</b>	<b>Southern-Ocean</b>
21	SAH	Sahara	42	EAU	E.Australia			

969



970 **Figures**

971  
972  
973  
974  
975  
976



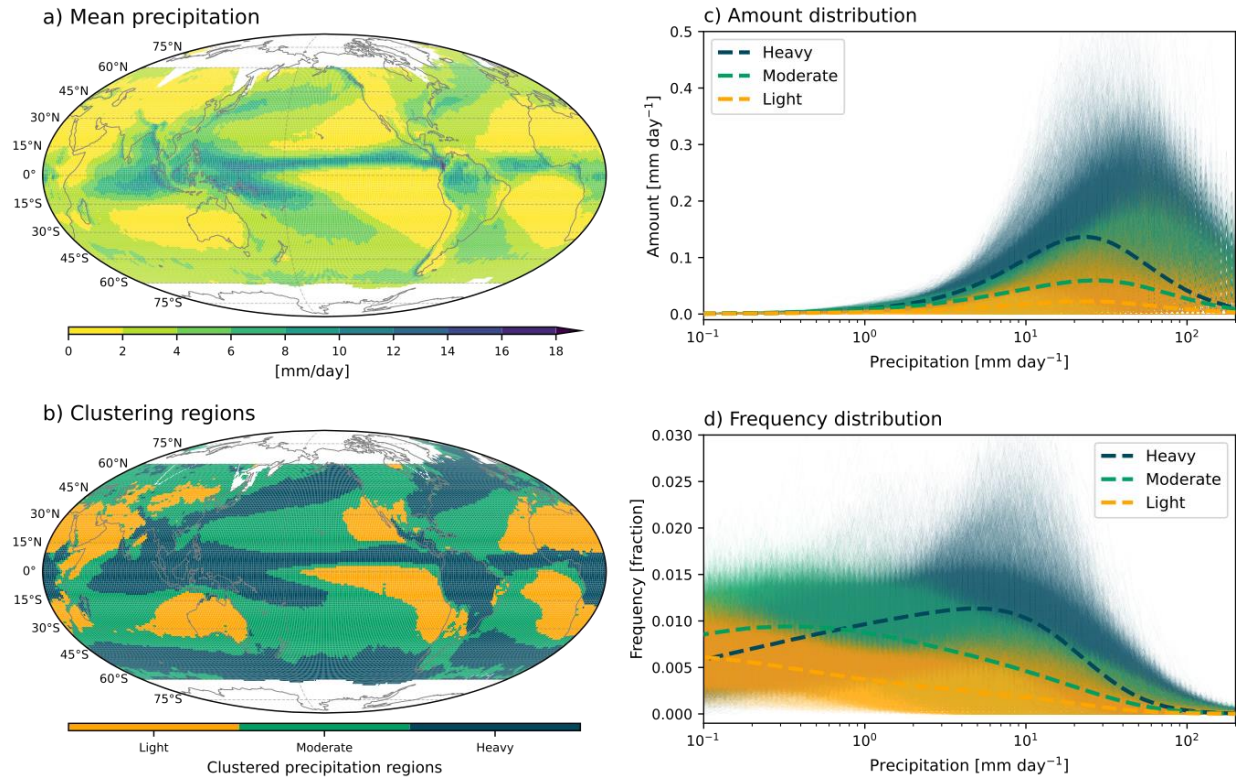
977  
978

979 Figure 1. Schematics for precipitation distribution metrics. a) Amount or Frequency  
980 distribution as a function of rain rate. Peak metric gauges the rain rate where the  
981 maximum distribution occurs. P10 and P90 metrics respectively measure the fraction of  
982 the distribution lower 10 percentile and upper 90 percentile. Perkins score is another  
983 metric based on the frequency distribution to quantify the similarity between observed  
984 and modeled distribution. b) Fraction of cumulative distribution as a function of number of  
985 wettest days. Unevenness gauges the number of wettest days for half of annual  
986 precipitation. FracPRdays measures the fraction of the number of precipitating  
987 ( $\geq 1\text{mm/day}$ ) days a year. SDII is designed to measure daily precipitation intensity by  
988 annual total precipitation divided by FracPRdays.

989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999

1000

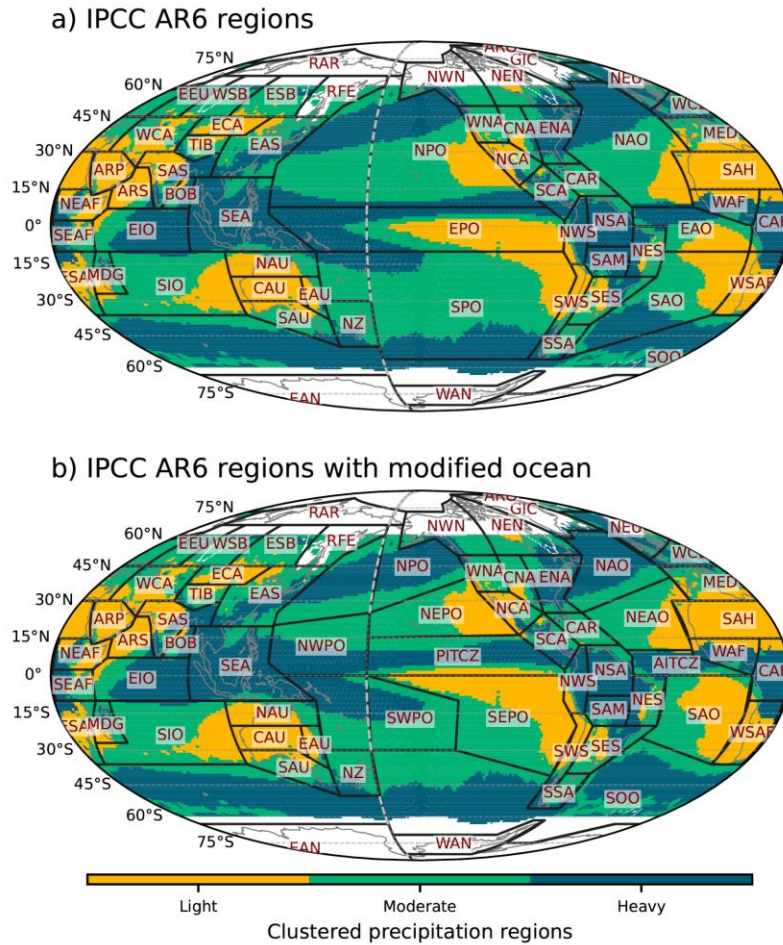
1001  
1002  
1003  
1004



1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023

Figure 2. Spatial patterns of IMERG precipitation a) mean state and b) clustering for heavy, moderate, and light precipitating regions by K-means clustering with amount and frequency distributions. Precipitation c) amount and d) frequency distributions as a function of rain rate. Different colors indicate different clustering regions as the same with b). Thin and thick curves respectively indicate distributions at each grid and the cluster average.

1024  
1025  
1026  
1027



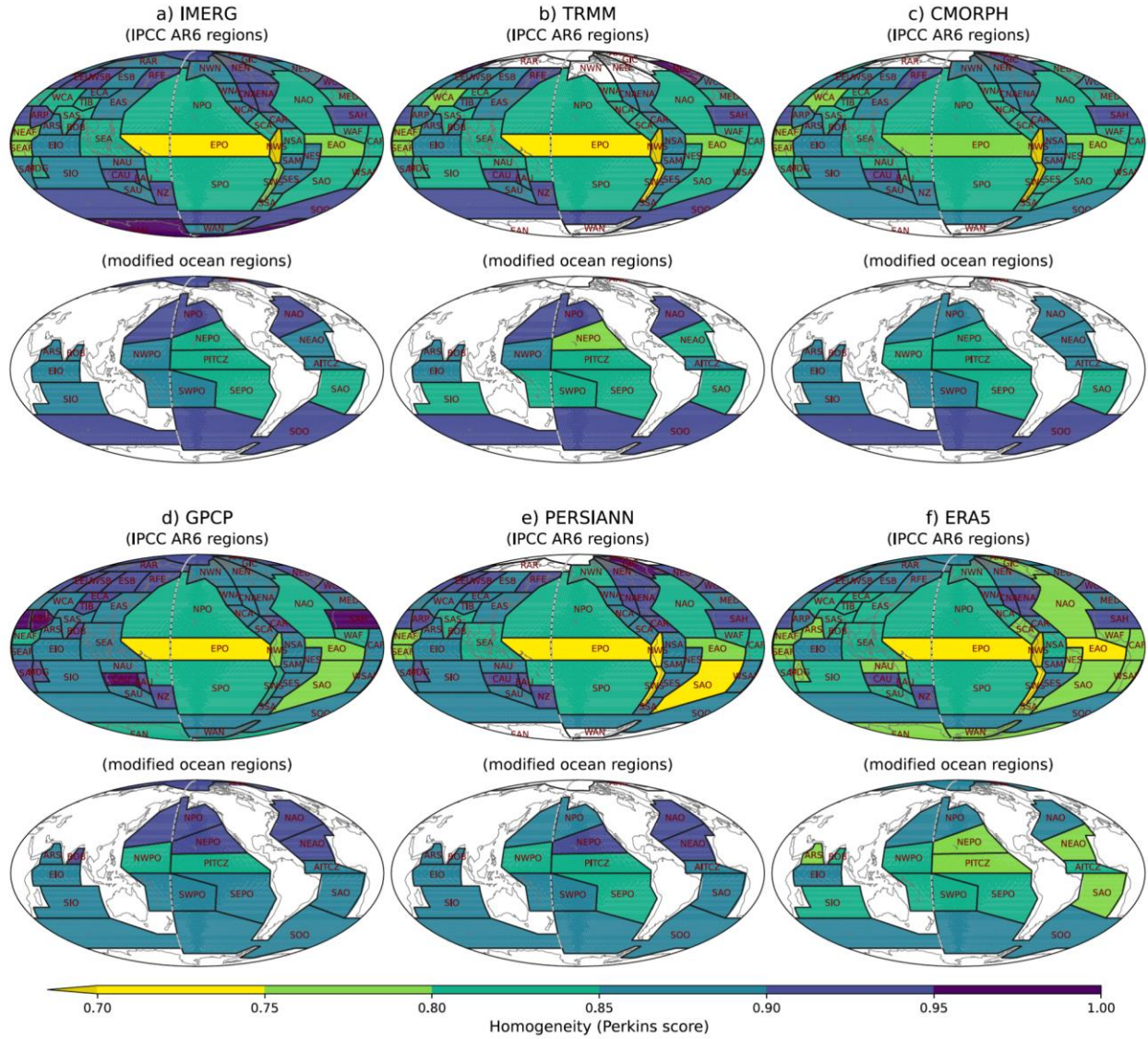
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042

Figure 3. a) IPCC AR6 climate reference regions and b) modified IPCC AR6 climate reference regions superimposed on the precipitation distributions clustering map shown in Fig. 2b. Land regions are the same between a) and b), while some ocean regions are modified.





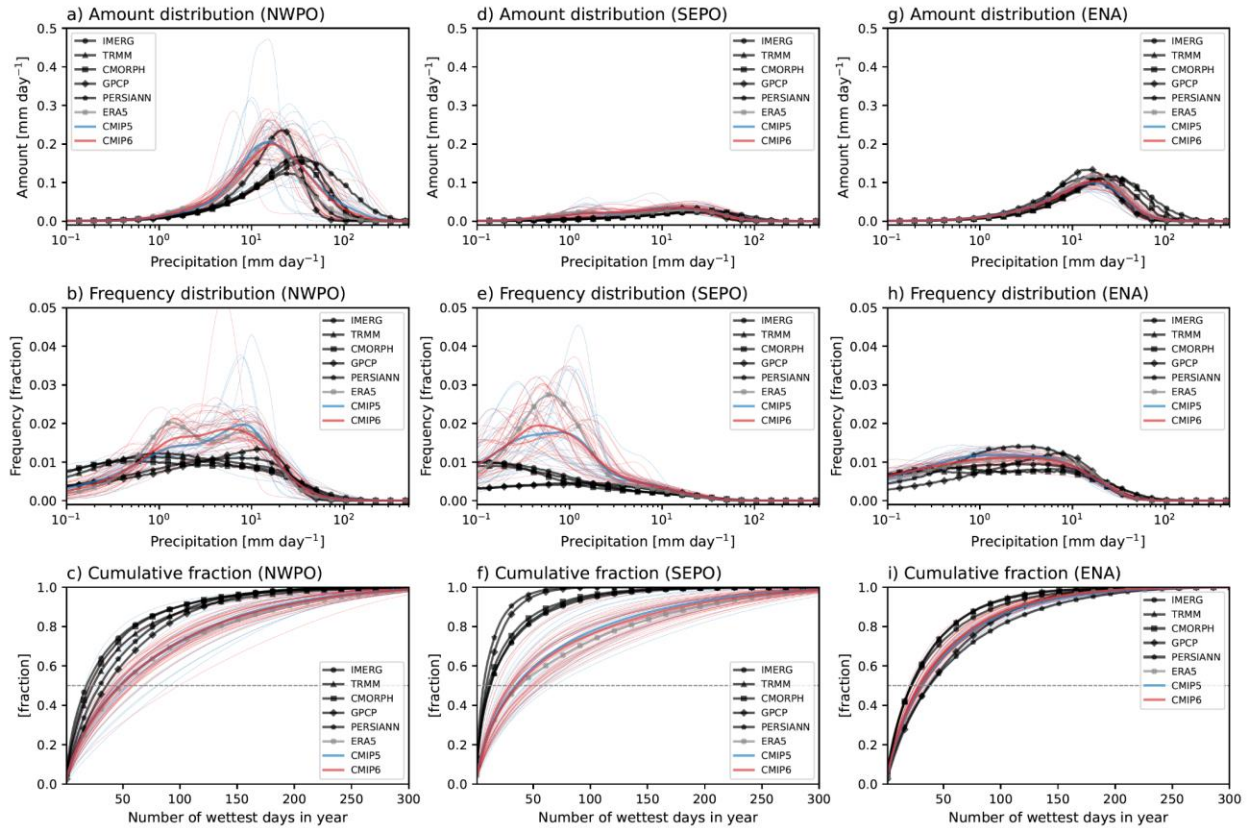
1071  
1072  
1073  
1074



1075  
1076  
1077  
1078  
1079  
1080  
1081  
1082  
1083  
1084  
1085

Figure 5. As in Fig. 4, but for different observational datasets with Perkins score.

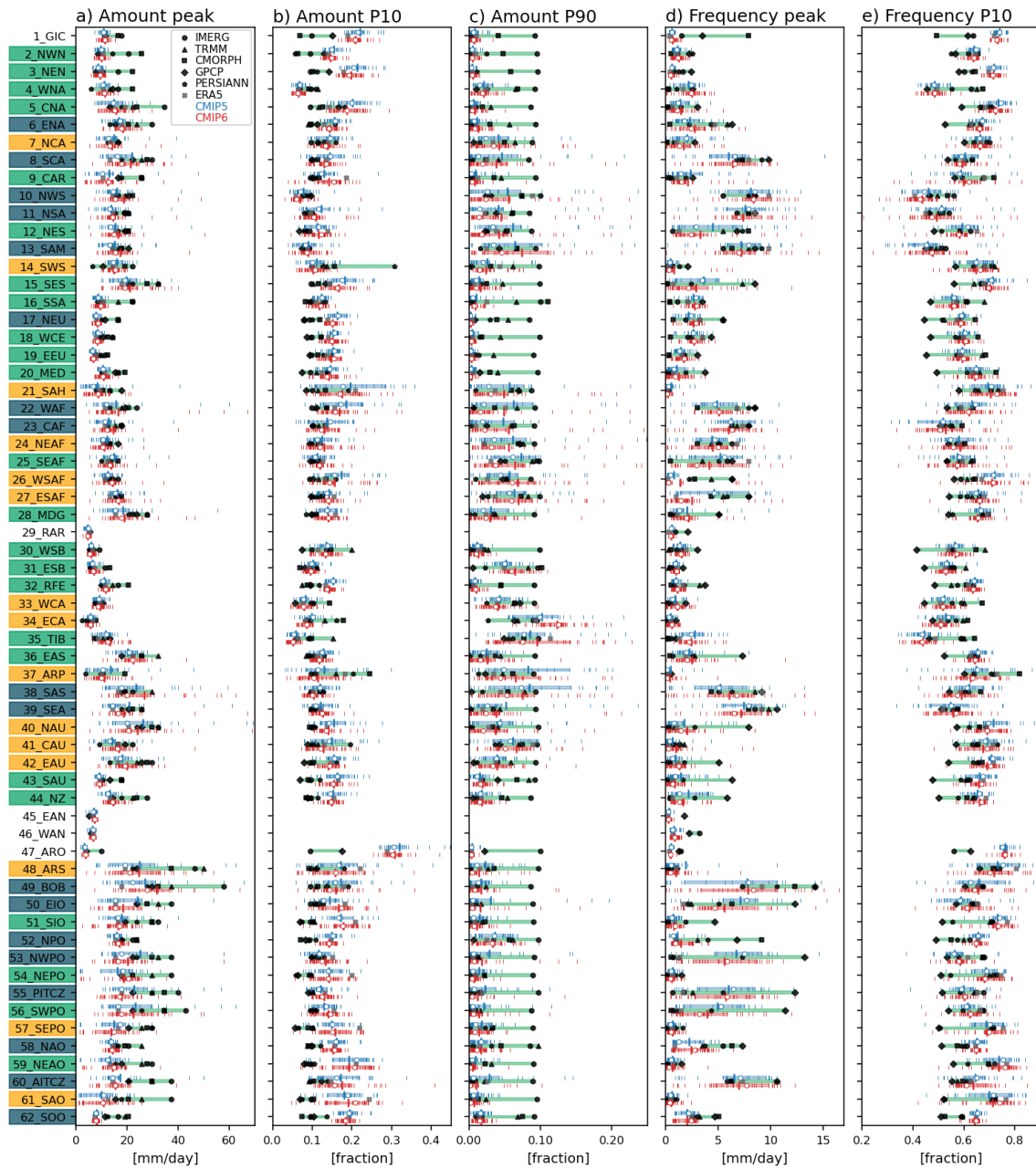
1086  
1087  
1088  
1089



1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107

Figure 6. Precipitation amount (upper), frequency (middle), and cumulative (bottom) distributions for a-c) NWPO, b-f) SEPO, and g-j) ENA. Black, gray, blue, and red curves indicate the satellite-based observations, reanalysis, CMIP5 models, and CMIP6 modes, respectively. Thin and thick curves for CMIP models respectively indicate distributions for each model and multi-model average. Gray dotted lines in the cumulative distributions indicate a fraction of 0.5. Note: all model output and observations were conservatively regridded to 2° in the first step of analysis.

1108  
1109  
1110

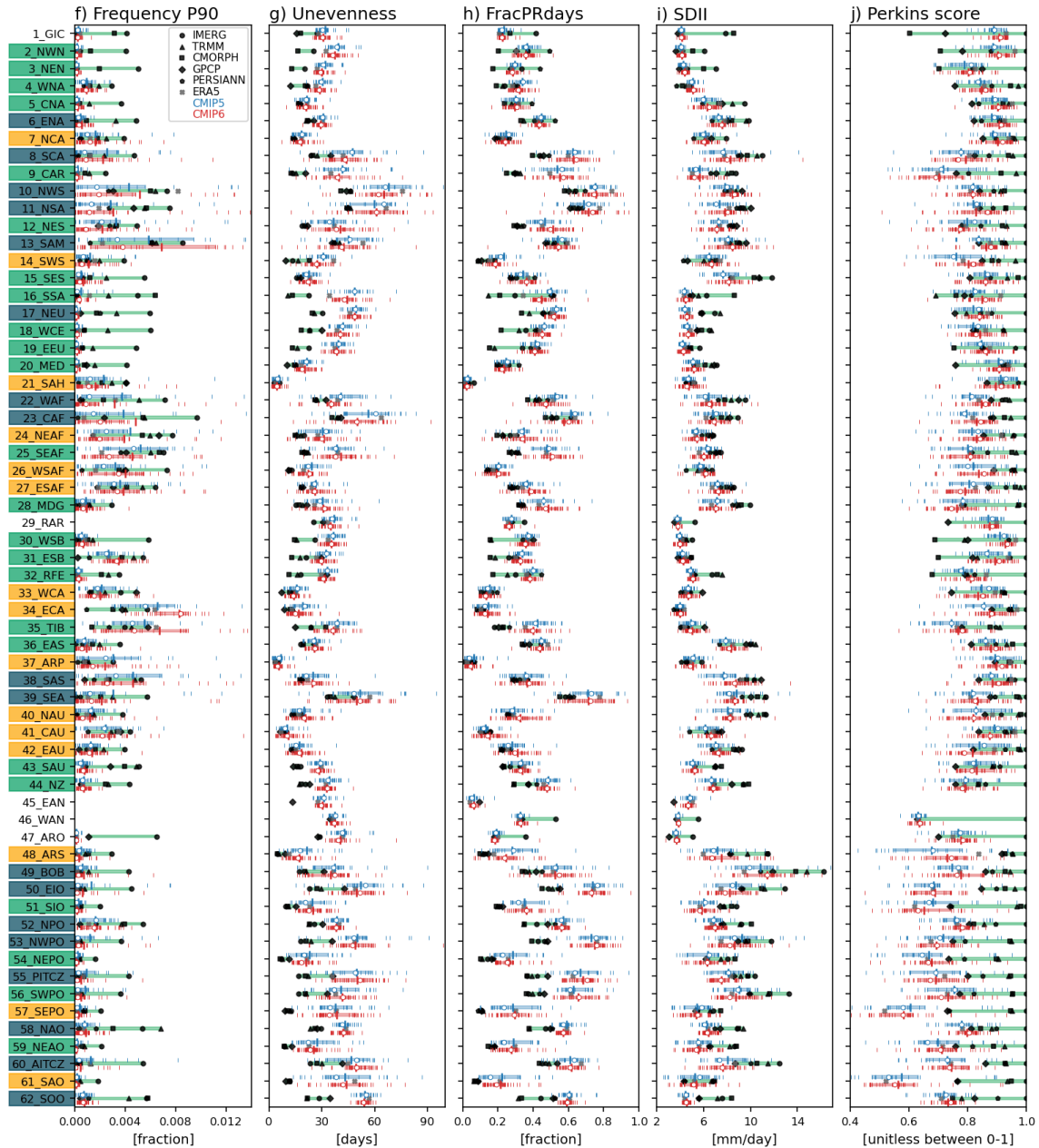


1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118

Figure 7. Precipitation distribution metrics for a) Amount peak, b) Amount P10, c) Amount P90, d) Frequency peak, e) Frequency P10, f) Frequency P90, g) Unevenness, h) FracPRdays, i) SDII, and j) Perkins score over the modified IPCC AR6 regions. Black, gray, blue, and red markers indicate the satellite-based observations, reanalysis, CMIP5 models, and CMIP6 modes, respectively. Thin and thick vertical marks for CMIP models respectively indicate distributions for each model and multi-model average.



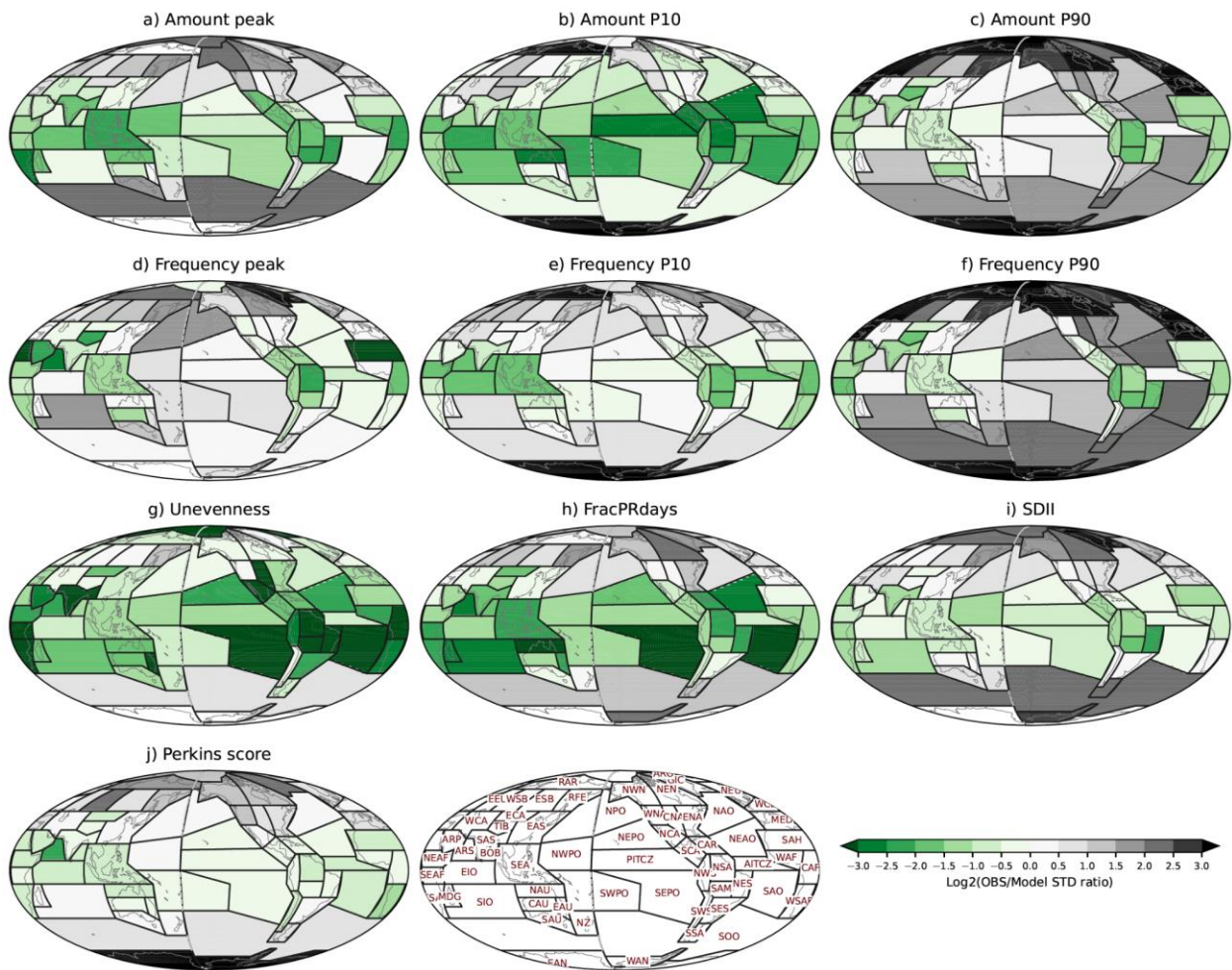
1119 Open circle mark for CMIP models indicates the multi-model median. Green shade  
 1120 represents the range between the minimum and maximum values of satellite-based  
 1121 observations. Blue and red shades respectively represent the range between 25th and  
 1122 75th model values for CMIP 5 and 6 models. Y-axis labels are shaded with the three  
 1123 colors as the same in Fig. 2b, indicating dominant precipitating characteristics. Note that  
 1124 regions 1-46 are land and land-ocean mixed regions, and 47-62 are ocean regions.  
 1125  
 1126



1127  
 1128  
 1129 Figure 7. (continued)



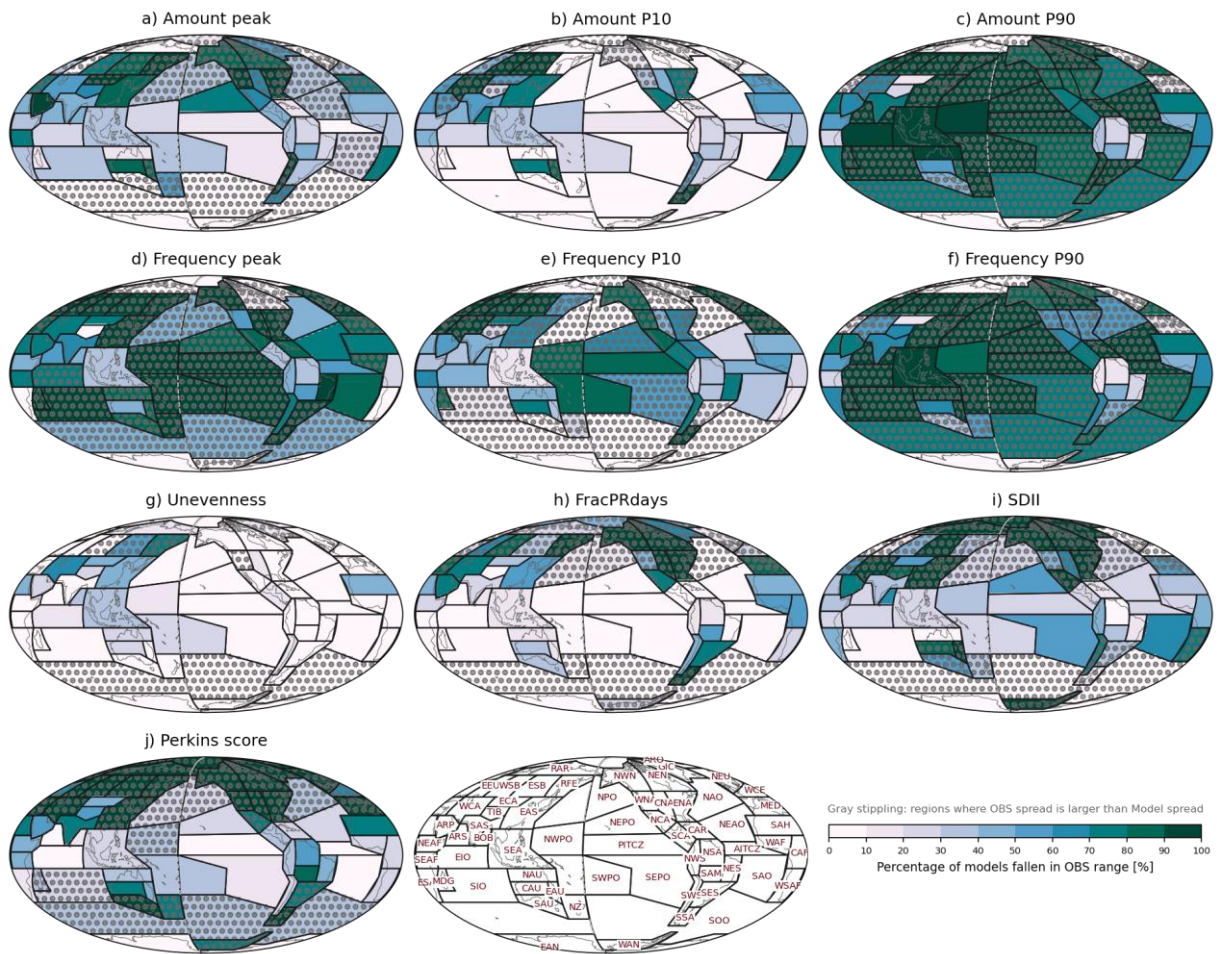
1130  
1131  
1132



1133  
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147

Figure 8. Observational discrepancies relative to spread in the multi-model ensemble for a) Amount peak, b) Amount P10, c) Amount P90, d) Frequency peak, e) Frequency P10, f) Frequency P90, g) Unevenness, h) FracPRdays, i) SDII, and j) Perkins score over the modified IPCC AR6 regions. The observational discrepancy is calculated by the standard deviation of satellite-based observations divided by the standard deviation of CMIP 5 and 6 models for each metric and region.

1148  
1149  
1150

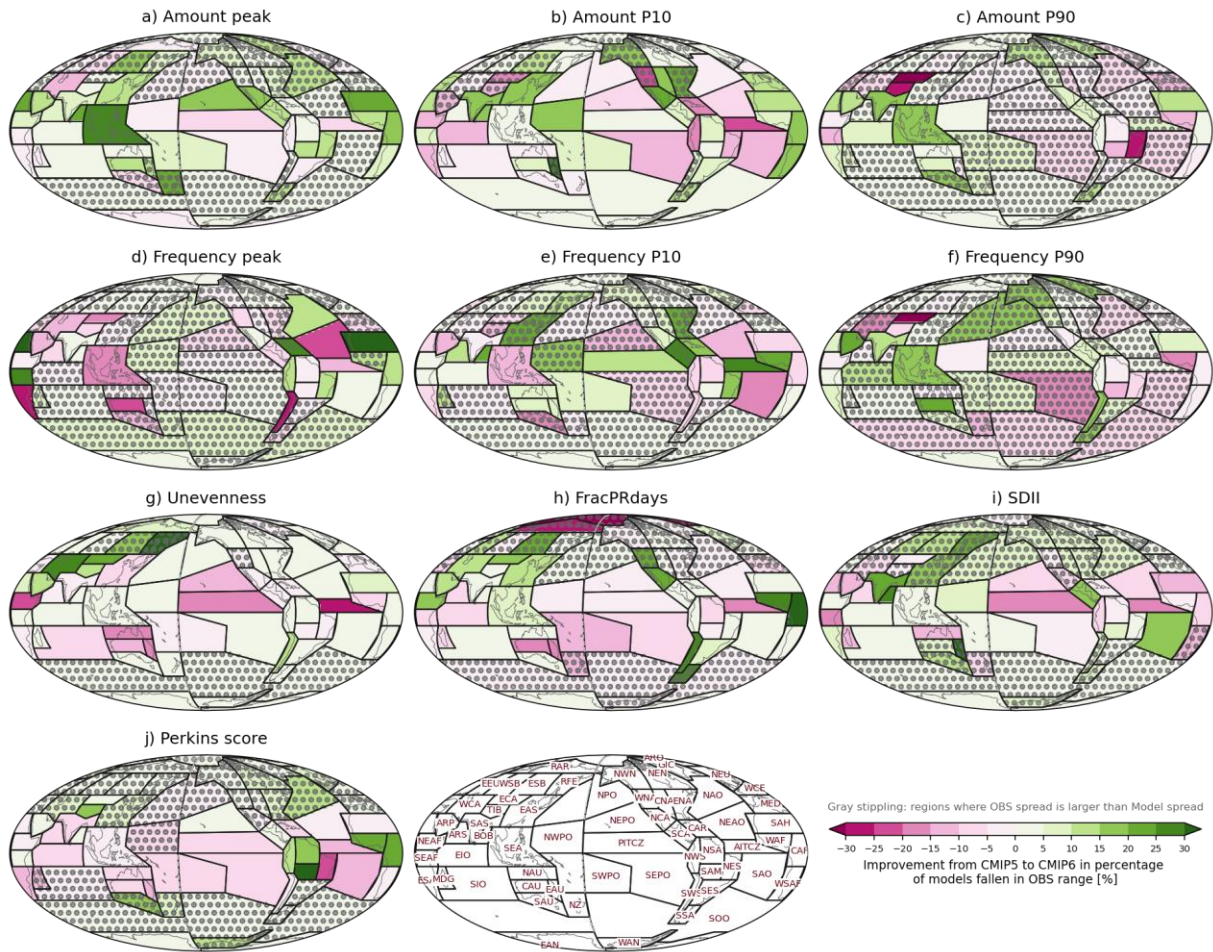


1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166

Figure 9. Percentage of CMIP6 models within range of the observational products for a) Amount peak, b) Amount P10, c) Amount P90, d) Frequency peak, e) Frequency P10, f) Frequency P90, g) Unevenness, h) FracPRdays, i) SDII, and j) Perkins score over the modified IPCC AR6 regions. The observational range is between the minimum and maximum values of five satellite-based products. Regions where the observational spread is larger than model spread shown in Fig. 8 are stippled gray.



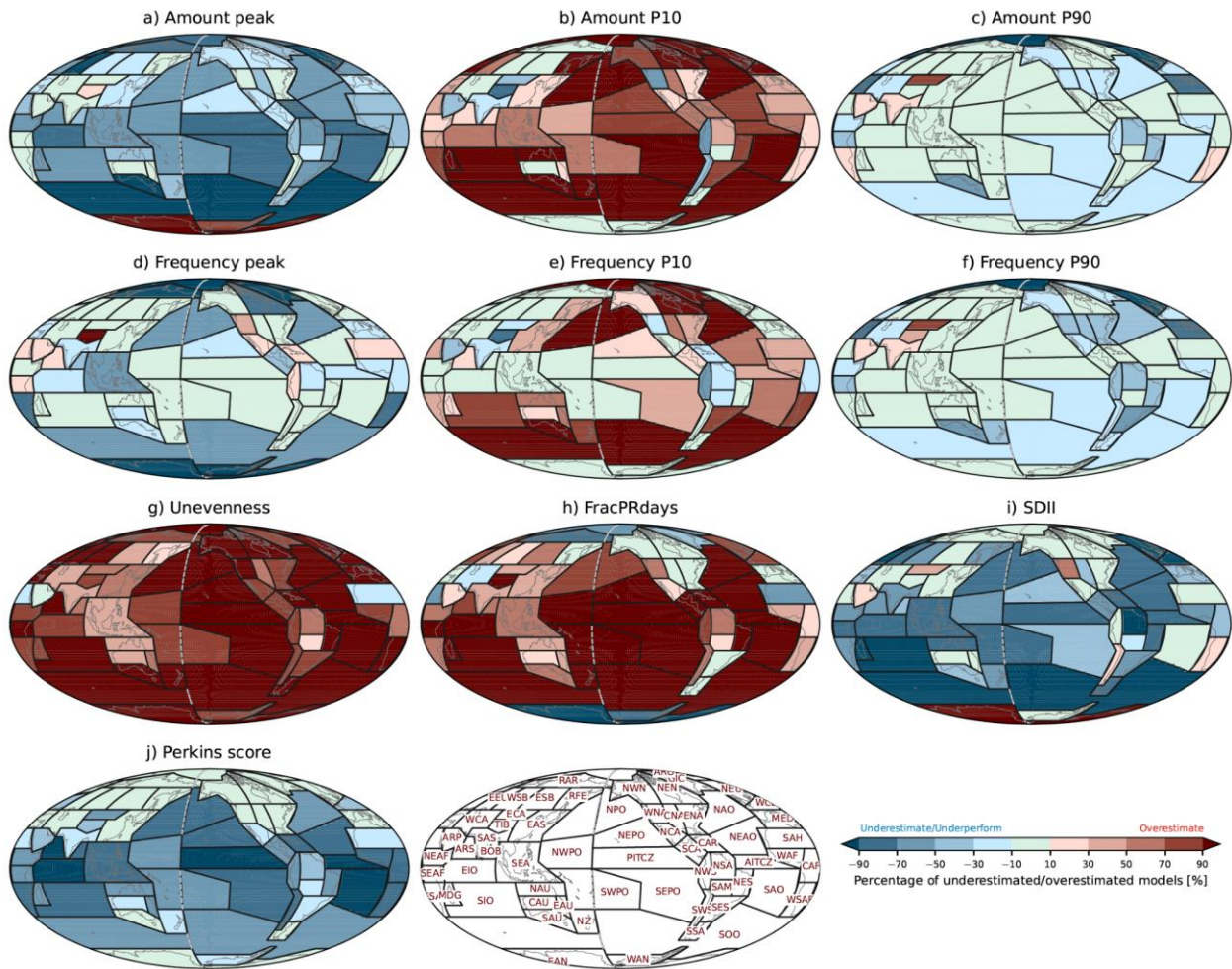
1167  
1168  
1169



1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185

Figure 10. Improvement from CMIP 5 to 6 as identified by the percentage of models in each multi-model ensemble that are within the observational min-to-max range. The improvement is calculated by the CMIP6 percentage minus the CMIP5 percentage, so that positive and negative values respectively indicate improvement and deterioration in CMIP6. Regions where the observational spread is larger than model spread are stippled gray.

1186  
1187  
1188

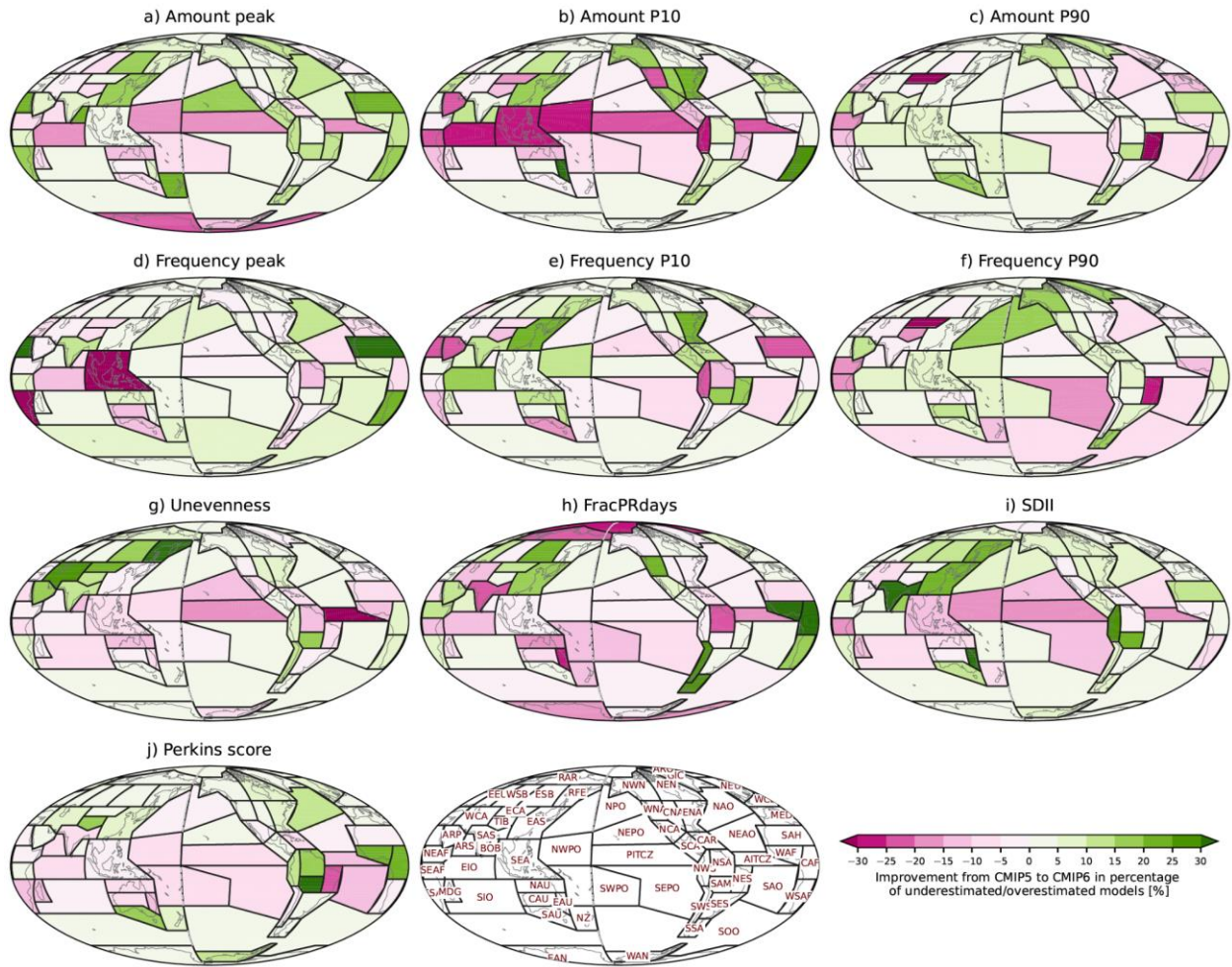


1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203

Figure 11. Percentage of CMIP6 models underestimating or overestimating observations for a) Amount peak, b) Amount P10, c) Amount P90, d) Frequency peak, e) Frequency P10, f) Frequency P90, g) Unevenness, h) FracPRdays, i) SDII, and j) Perkins score over the modified IPCC AR6 regions. The criteria for underestimation and overestimation are respectively defined by minimum and maximum values of satellite-based observations shown in Fig. 7. Positive and negative values respectively represent overestimation and underestimation by a formulation of  $(nO - nU)/nT$  where  $nO$ ,  $nU$ ,  $nT$  are respectively the number of overestimated models, underestimated models, and total models.



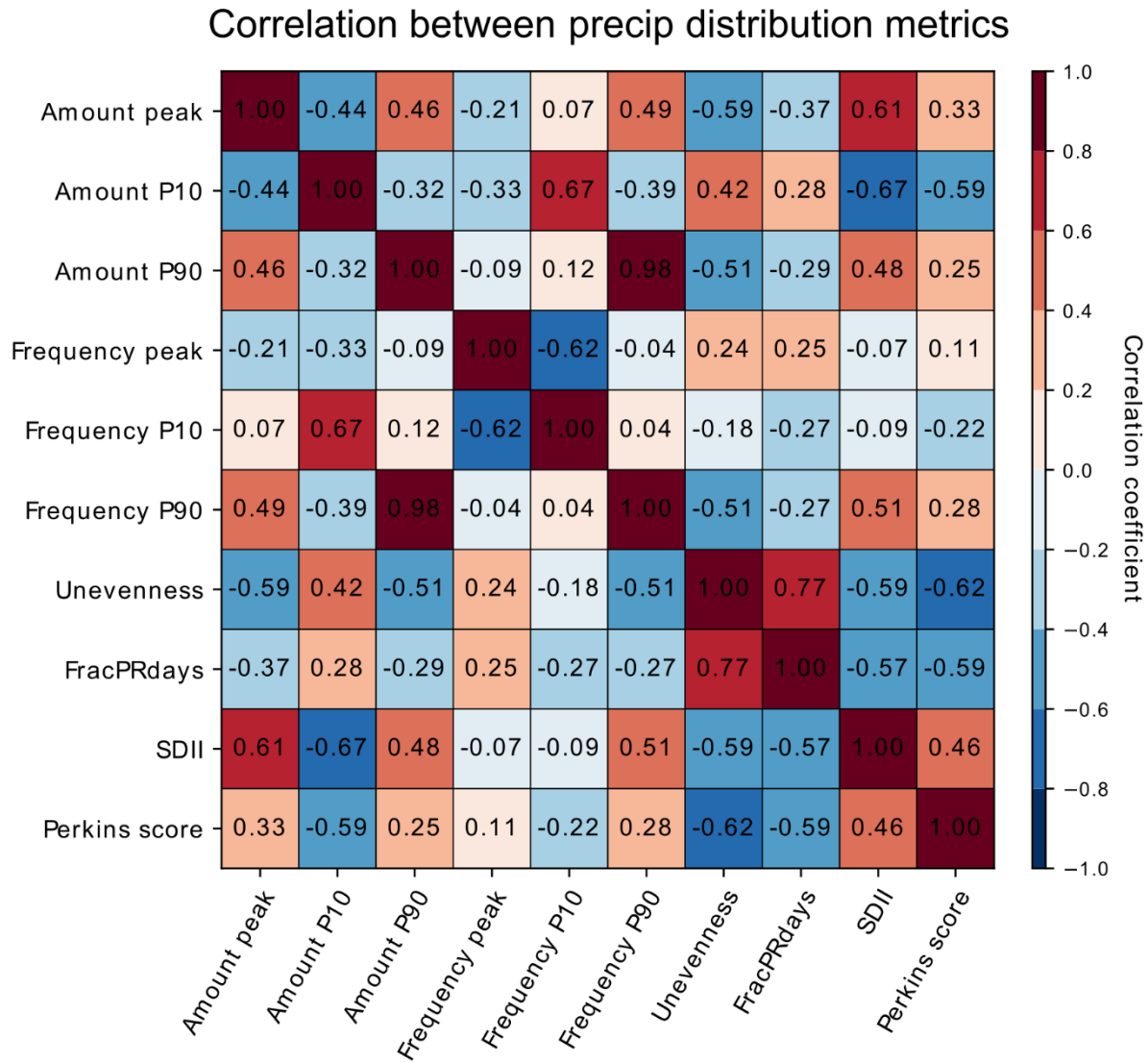
1204  
1205  
1206



1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221

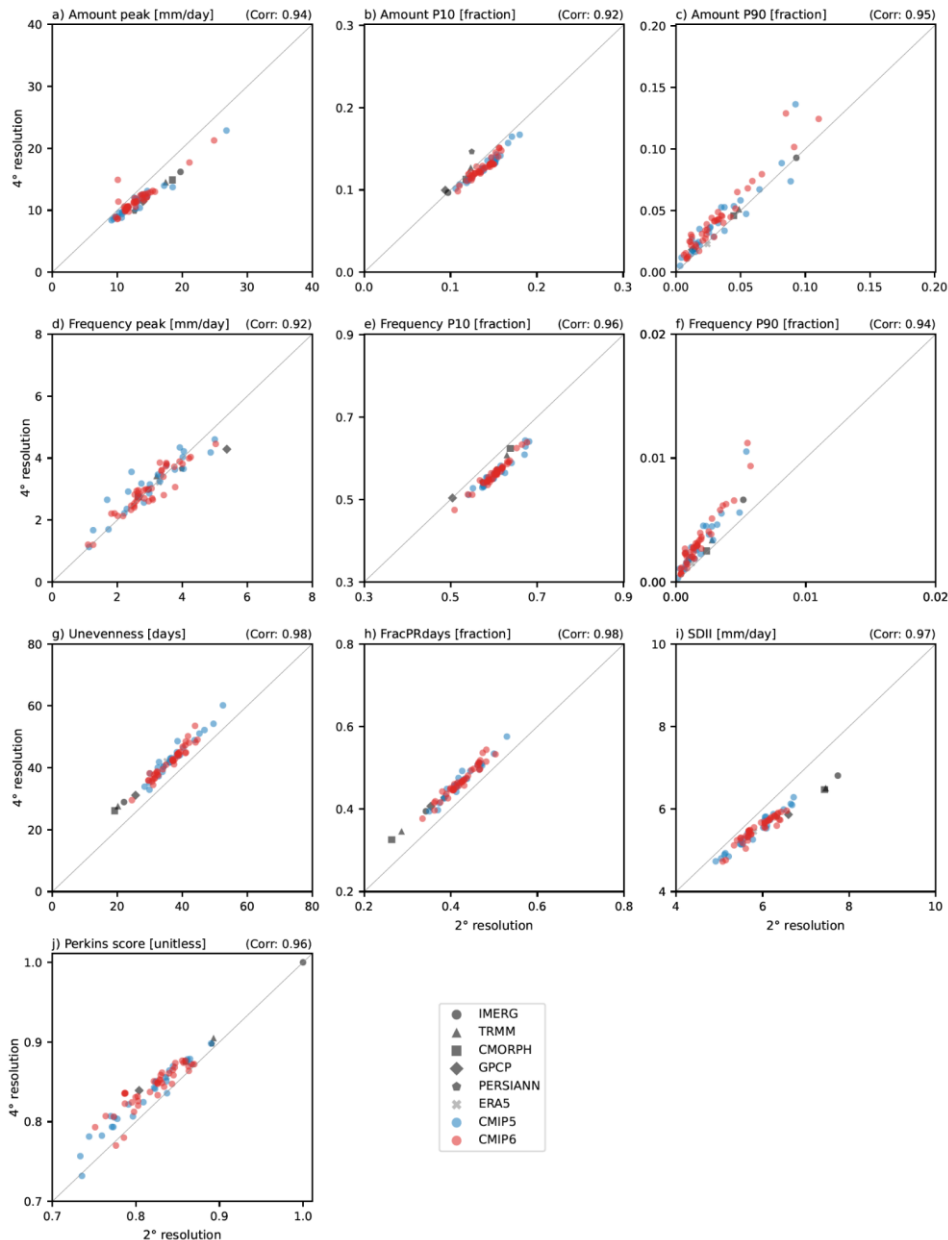
Figure 12. Improvement from CMIP 5 to 6 in the percentage of underestimated or overestimated models. The improvement is calculated by the absolute value of CMIP5 percentage minus the absolute value of CMIP6 percentage, so that positive and negative values respectively indicate improvement and deterioration in CMIP6.

1222  
 1223  
 1224



1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236

Figure 13. Correlation between precipitation distribution metrics across CMIP 5 and 6 model performances. The correlation coefficients are calculated for the modified IPCC AR6 regions and then area-weighted averaged globally.



1238

1239

1240 Figure 14. Scatterplot between 2° and 4° interpolated horizontal resolutions in  
 1241 evaluating precipitation distribution metrics for a) Amount peak, b) Amount P10, c)  
 1242 Amount P90, d) Frequency peak, e) Frequency P10, f) Frequency P90, g) Unevenness,  
 1243 h) FracPRdays, i) SDII, and j) Perkins score. The metric values are calculated for the  
 1244 modified IPCC AR6 regions and then weighted averaged globally. Black, gray, blue, and  
 1245 red marks indicate the satellite-based observations, reanalysis, CMIP5 models, and  
 1246 CMIP6 modes, respectively. The number in the upper right of each panel is the  
 1247 correlation coefficient between the metric values in 2° and 4° resolutions across all  
 1248 observations and models.

1249  
1250