



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

Min-Seop Ahn^{1,*}, Paul A. Ullrich^{2,1}, Peter J. Gleckler¹, Jiwoo Lee¹, Ana C. Ordonez¹, and Angeline G. Pendergrass^{3,4}

¹PCMDI, Lawrence Livermore National Laboratory, Livermore, CA, USA

²Department of Land, Air and Water Resources, University of California, Davis, CA, USA

³Earth and Atmospheric Science, Cornell University, Ithaca, NY, USA

⁴National Center for Atmospheric Research, Boulder, CO, USA

* Corresponding author: Min-Seop Ahn (ahn6@llnl.gov)



1 **Abstract**

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

21



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. Pendergrass et al. (2020) envisioned a framework for both
44 baseline and exploratory precipitation benchmarks, and Leung et al. (2022) described



45 efforts to advance exploratory objective evaluation for simulated precipitation focusing on
46 process-oriented and phenomena-based metrics at a variety of spatiotemporal scales.
47 The baseline precipitation benchmark metrics target established measures of the mean
48 state, the seasonal and diurnal cycles, variability across timescales, intensity/frequency
49 distributions, extremes, and drought. The current study provides a framework focused on
50 precipitation distributions.

51

52 Diagnosing precipitation distributions and formulating metrics that extract critical
53 information from precipitation distributions have been addressed in many previous studies.
54 Pendergrass and Deser (2017) proposed several precipitation distribution metrics based
55 on frequency and amount distribution curves. The precipitation frequency distribution
56 quantifies how often rain occurs at different rain rates, whereas the precipitation amount
57 distribution quantifies how much rain falls at different rain rates. Based on the distribution
58 curves, Pendergrass and Deser (2017) extracted rain frequency peak and amount peak
59 where the maximum non-zero rain frequency and amount occur, respectively.
60 Pendergrass and Knutti (2018) introduced a metric that measures the unevenness of daily
61 precipitation based on the cumulative amount curve. Their unevenness metric is defined
62 as the number of wettest days that constitute half of the annual precipitation. In the
63 median of station observations equatorward of 50° latitude, half of the annual precipitation
64 falls in only about the heaviest 12 days, and generally models underestimate the
65 observed unevenness (Pendergrass and Knutti 2018). In addition, several metrics have
66 been developed to distill important precipitation characteristics, such as the fraction of
67 precipitating days and simple daily intensity index (SDII, Zhang et al. 2011). In this study



68 we implement all these well-established metrics and several other complementary metrics
69 into our framework.

70

71 Many studies have analyzed the precipitation distributions over large domains (e.g., Dai
72 2006; Pendergrass and Hartmann 2014; Ma et al. 2022). Often, these domains comprise
73 both heavily precipitating and dry regions. Given the emphasis on regional scale analysis
74 continues to grow as models' horizontal resolution increases, interpretation of domain-
75 averaged distributions could be simplified by defining regions that are not overly complex
76 or heterogeneous in terms of their precipitation distribution characteristics. Iturbide et al.
77 (2020) has identified climate reference regions that have been adopted in the sixth
78 assessment report (AR6) of the Intergovernmental Panel on Climate Change (IPCC). Our
79 framework is based on these IPCC AR6 reference regions for objective examination of
80 precipitation distributions over land. Over the ocean we have revised some of the regions
81 of Iturbide et al. (2020) to better isolate homogeneous precipitation distribution
82 characteristics.

83

84 In this study, we propose a framework for regional scale quantification of simulated
85 precipitation distributions and evaluate CMIP 5 and 6 models with multiple observations.
86 The remainder of the paper is organized as follows: Sections 2 and 3 describe the data
87 and analysis methods. Section 4 presents results including the application and
88 modification of IPCC AR6 climate reference regions, evaluation of CMIP models, and
89 their improvement across generations. Sections 5 and 6 discuss and summarize the main
90 accomplishments and findings from this study.



91

92

93 **2. Data**

94 2.1. Observational data

95 For reference data, we use six global daily precipitation products first made available as
96 part of the Frequent Rainfall Observations on GridS (FROGS) database (Roca et al., 2019)
97 and then further aligned with CMIP output via the data specifications of the Observations
98 for Model Intercomparison Project (Obs4MIPs, Waliser et al. 2020). These include five
99 satellite-based products and a recent atmospheric reanalysis product. The satellite-based
100 precipitation products include the Integrated Multi-satellitE Retrievals for GPM version 6
101 final run product (Huffman et al. 2020; hereafter IMERG), the Tropical Rainfall Measuring
102 Mission Multi-satellite Precipitation Analysis 3B42 version 7 product (Huffman et al. 2007;
103 hereafter TRMM), the bias-corrected Climate Prediction Center Morphing technique
104 product (Xie et al. 2017; hereafter CMORPH), the Global Precipitation Climatology Project
105 1DD version 1.3 (Huffman et al. 2001; hereafter GPCP), and Precipitation Estimation from
106 Remotely Sensed Information using Artificial Neural Networks (Ashouri et al. 2015;
107 hereafter PERSIANN). The reanalysis product included for context is the ECMWF's fifth
108 generation of atmospheric reanalysis (Hersbach et al. 2020; hereafter ERA5). Table 1
109 summarizes the observational datasets with the data source, coverage of domain and
110 period, resolution of horizontal space and time frequency, and references. We use the
111 data periods available via FROGS and Obs4MIPs as follows: 2001-2020 for IMERG,
112 1998-2018 for TRMM, 1998-2012 for CMORPH, 1997-2020 for GPCP, 1984-2018 for
113 PERSIANN, and 1979-2018 for ERA5.



114

115 2.2. CMIP model simulations

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

125

126

127 3. Methods

128 In our framework we apply 10 metrics that characterize different and complementary
129 aspects of the intensity distribution of precipitation at regional scales. Table 3 summarizes
130 the metrics including their definition, purpose, and references. The computation of the
131 metrics has been implemented and applied in an open source metrics package, the
132 Program for Climate Model Diagnosis & Intercomparison (PCMDI) metrics package (PMP;
133 Gleckler et al. 2008, 2016).

134

135 3.1. Frequency and amount distributions



136 Following Pendergrass and Hartmann (2014) and Pendergrass and Deser (2017), we use
137 logarithmically-spaced bins of daily precipitation to calculate both the precipitation
138 frequency and amount distributions. Each bin is 7% wider than the previous one, and the
139 smallest non-zero bin is centered at 0.03 mm/day. The frequency distribution is the
140 number of days in each bin normalized by the total number of days, and the amount
141 distribution is the sum of accumulated precipitation in each bin normalized by the total
142 number of days. Based on these distributions (Fig. 1a), we identify the rain rate where the
143 maximum peak of the distribution appears (Amount/Frequency Peak, Pendergrass and
144 Deser 2017; also called mode, Kooperman et al., 2016) and formulate several
145 complementary metrics that measure the fraction of the distribution lower 10 percentile
146 (P10) and upper 90 percentile (P90). The precipitation bins less than 0.1 mm/day are
147 considered dry for the purpose of these calculations. The threshold rain rates for 10th and
148 90th percentiles are defined from the amount distribution as determined from
149 observations. Here we use IMERG as the default reference observational dataset. The
150 final frequency related metric we employ is the Perkins score, which measures the
151 similarity between observed and modeled frequency distributions (Perkins et al. 2007).
152 With the sum of a frequency distribution across all bins being unity, the Perkins score is
153 defined as the sum of minimum values between observed and modeled frequency across
154 all bins: $Perkins\ Score = \sum_1^n \text{minimum}(Z_o, Z_m)$ where n is the number of bins, Z_o and
155 Z_m are the frequency in a given bin for observation and model, respectively. The Perkins
156 score is a unitless scalar varying from 0 (low similarity) to 1 (high similarity).

157

158 3.2. Cumulative fraction of annual precipitation amount



159 Following Pendergrass and Knutti (2018), we calculate the cumulative sum of daily
160 precipitation each year sorted in descending order (i.e., wettest to driest) and normalized
161 by the total precipitation for that year. From the distribution for each individual year (see
162 Fig. 1b), we obtain the metrics gauging the number of the wettest days for half of annual
163 precipitation (Unevenness, Pendergrass and Knutti 2018) and the fraction of the number
164 of precipitating (≥ 1 mm/day) days (FracPRdays). To facilitate comparison against longer-
165 established analyses (e.g., ETCCDI, Zhang et al., 2011), we include the daily
166 precipitation intensity, calculated by dividing the annual total precipitation by the number
167 of precipitating days (SDII, Zhang et al. 2011). To obtain values of these metrics over
168 multiple years, we take the median across years following Pendergrass and Knutti (2018;
169 for unevenness).

170

171 3.3. Reference regions

172 We use the spatial homogeneity of precipitation characteristics as a basis for defining
173 regions, as in previous studies (e.g., Swenson and Grotjahn 2019). In addition to
174 providing more physically-based results, this also simplifies interpretation with robust
175 diagnostics when we average a distribution characteristic across the region. We use K-
176 means clustering (MacQueen 1967) with the concatenated frequency and amount
177 distributions of IMERG over the global domain to identify homogeneous regions for
178 precipitation distributions. K-means clustering is an unsupervised machine learning
179 algorithm that separates characteristics of a given dataset into a specified number of
180 groups, which has been widely used because it is faster and simpler than other methods.
181 Figure 2 shows the spatial pattern of IMERG precipitation mean state and clustering



182 results with 3 clusters identified by the algorithm (Fig. 2b) including heavy (blue),
183 moderate (green), and light (orange) precipitation regions. The spatial pattern of these
184 clustering regions resembles the pattern of the mean state of precipitation, providing a
185 sanity check indicating that the cluster-based regions are physically reasonable.

186

187 In support of the AR6, the IPCC proposed a set of climate reference regions (Iturbide et
188 al. 2020). These regions were defined based on geographical and political boundaries
189 and the climatic consistency of temperature and precipitation in current climate and
190 climate change projections. When defining regions, the land regions use both information
191 from current climate and climate change projections, while the ocean regions use only
192 the information from climate change projections. In other words, the climatic consistency
193 of precipitation in the current climate is not explicitly represented in the definition of the
194 oceanic regions. Figure 3a shows the IPCC AR6 climate reference regions superimposed
195 on our precipitation clustering regions shown in Fig. 2b. The land regions correspond
196 reasonably well to the clustering regions, but some ocean regions are too broad, including
197 both heavy and light precipitating regions (Fig. 3a). In this study, the ocean regions are
198 modified based on the clustering regions, while the land regions remain the same as in
199 the AR6 (Fig. 3b).

200

201 In the Pacific Ocean region, the original IPCC AR6 regions consist of equatorial Pacific
202 Ocean (EPO), northern Pacific Ocean (NPO), and southern Pacific Ocean (SPO). Each
203 of these regions includes areas of both heavy and light precipitation. EPO includes the
204 Intertropical Convergence Zone (ITCZ), the South Pacific Convergence Zone (SPCZ),



205 and also the dry southeast Pacific region. The NPO region includes the north Pacific storm
206 track and the dry northeast Pacific. The SPO region includes the southern part of SPCZ
207 and the dry southeast area of the Pacific. In our modified IPCC AR6 regions, the Pacific
208 Ocean region is divided into four heavy precipitating regions (NPO, NWPO, PITCZ, and
209 SWPO) and two light and moderate precipitating regions (NEPO and SEPO). The NPO,
210 NWPO, PITCZ, and SWPO mainly include the North Pacific storm track region, the
211 western Pacific warm pool region, pacific ITCZ, and SPCZ, respectively. The NEPO and
212 SEPO respectively include the northeast and southeast dry Pacific regions. Similarly, in
213 the Atlantic Ocean region, the original IPCC AR6 regions consist of the equatorial Atlantic
214 Ocean (EAO), northern Atlantic Ocean (NAO), and southern Atlantic Ocean (SAO), with
215 each including both heavy and light precipitating regions. Our modified Atlantic Ocean
216 region consists of two heavy precipitating regions (NAO and AITCZ) and two light and
217 moderate precipitating regions (NEAO and SAO). The NAO and AITCZ mainly include
218 the North Atlantic storm track region and Atlantic ITCZ, respectively. The NEAO and SAO
219 mainly include dry eastern Atlantic regions. The Indian Ocean (IO) region is not modified
220 as the original IPCC AR6 climate reference region separates well the heavy precipitating
221 equatorial IO (EIO) region from the moderate and light precipitating southern IO (SIO)
222 region. The Southern Ocean (SOO) is modified to mainly include the heavy precipitation
223 region around the Antarctic. The original IPCC AR6 climate reference regions consist of
224 58 regions including 12 oceanic regions and 46 land regions, while our modification
225 consists of 62 regions including 16 oceanic regions and the same land regions as the
226 original (see Table 4). Note that the Caribbean (CAR), the Mediterranean (MED), and
227 Southeast Asia (SEA) are not counted for the oceanic regions.



228

229 3.4. Evaluating model performance

230 We use two simple measures to compare the collection of CMIP 5 and 6 model
231 simulations with the five satellite-based observational products (IMERG, TRMM,
232 CMORPH, GPCP, and PERSIANN). One gauges how many models within the multi-
233 model ensemble fall within the observational range between the minimum and maximum
234 observed values for each metric and each region. Another is how many models
235 underestimate or overestimate all observations, i.e., are outside the bounds spanned by
236 the minimum and maximum values across the five satellite-based products. To quantify
237 the dominance of underestimation versus overestimation of the multi-model ensemble
238 with a single number, we use the following measure formulation: $(nO - nU)/nT$ where nO
239 is the number of overestimating models, nU is the number of underestimating models,
240 and nT is the total number of models. Thus, positive values represent overestimation, and
241 negative values represent underestimation. If models are mostly within the observational
242 range or widely distributed from underestimation to overestimation, the quantification
243 value would approach zero.

244

245 Many metrics that can be used to characterize precipitation, including those used here,
246 are sensitive to the spatial and temporal resolutions at which the model and observational
247 data are analyzed (e.g., Pendergrass and Knutti 2018, Chen and Dai 2019). As in many
248 previous studies the diagnosis of precipitation in CMIP 5 and 6 models (e.g., Fiedler et al.
249 2020; Tang et al. 2021; Ahn et al. 2022), to ensure appropriate comparisons, we conduct
250 all analyses at a common horizontal grid of 2x2 degrees with a conservative regridding



251 method. For models with multiple ensemble members, we first compute the metrics for
252 all available realizations and then average the results across the realizations.

253

254

255 **4. Results**

256 4.1. Homogeneity within reference regions

257 For the regional scale analysis, we employ the IPCC AR6 climate reference regions
258 (Iturbide et al. 2020) while we revise the region dividings over the oceans based on
259 clustered precipitation characteristics as described in section 3.3. To quantitatively
260 evaluate the homogeneity of precipitating distributions in the reference regions, we use
261 three homogeneity metrics: the Perkins score (Perkins et al. 2007), Kolmogorov–Smirnov
262 test (K-S test, Chakravart et al. 1967), and Anderson-Darling test (A-D test, Stephens
263 1974). The three metrics measure the similarity between the regionally-averaged and
264 individual grid cell frequency distributions within the region. The Perkins score measures
265 the overall similarity between two frequency distributions, which is one of our distribution
266 performance metrics described in Section 3.1. The K-S and A-D tests focus more on the
267 similarity in the center and the side of the frequency distribution, respectively. The three
268 homogeneity metrics could complement each other as their main focuses are on different
269 aspects of frequency distributions.

270

271 In the original IPCC AR6 reference regions, the oceanic regions show relatively low
272 homogeneity of precipitating characteristics compared to land regions (Fig. 4). The Pacific
273 and Atlantic Ocean regions show much lower homogeneity than the Indian Ocean,



274 especially in EPO and EAO regions. In the modified oceanic regions, the homogeneities
275 show an overall improvement with the three homogeneity metrics. In particular, the
276 homogeneity over the heavy precipitating regions where the homogeneity was lower (e.g.,
277 Pacific and Atlantic ITCZ and mid-latitude storm track regions) are largely improved. The
278 clustering regions shown here are obtained based on IMERG precipitation. However,
279 since different satellite-based products show substantial discrepancies in precipitation
280 distributions, it is important to assess whether the improved homogeneity in the modified
281 regions is similarly improved across other different datasets. Figure 5 shows the
282 homogeneity of precipitation distribution characteristics for different observational
283 datasets using the Perkins score. Although the region modifications we have made are
284 based on the clustering regions of IMERG precipitation, Fig. 5 suggests that the
285 improvement of the homogeneity over the modified regions is consistent across different
286 observational datasets. We further tested the homogeneity for different seasons (see Fig.
287 S1 in the supplement material). The homogeneity is overall improved in the modified
288 regions across the seasons even though we defined the reference regions based on
289 annual data.

290

291 4.2. Regional evaluation of model simulations against multiple observations

292 The precipitation distribution metrics used in this study are mainly calculated from three
293 curves: amount distribution, frequency distribution, and cumulative amount fraction
294 curves. Figure 6 shows these curves for three selected regions based on the clustered
295 precipitating characteristics (NWPO, which is a heavy precipitation dominated ocean
296 region; SEPO, a light precipitation dominated ocean region; and ENA, a heavy



297 precipitation dominated land region). The heavy and light precipitating regions are well
298 distinguished by their overlaid distribution curves. The amount distribution has a
299 distinctive peak in the heavy precipitating region (Figs. 6a and 6g), while it is flatter in the
300 light precipitating region (Fig. 6d). The frequency distribution is more centered on the
301 heavier precipitation side in the heavy precipitating region (Figs. 6b, 6h) than in the light
302 precipitating region (Fig 6e). The cumulative fraction increases more steeply in the light
303 precipitating region (Fig. 6f) than in the heavy precipitating region (Figs. 6c and 6i),
304 indicating there are fewer precipitating days in the light precipitating region. NWPO and
305 SEPO were commonly averaged for representing the tropical ocean region in many
306 studies, but these different characteristics in the precipitation distributions demonstrate
307 the additional information available via a regional scale analysis. Although satellite-based
308 observations are less reliable over the light precipitating ocean regions (e.g., SEPO), the
309 differences between heavy and light precipitation regions are well distinguishable.

310

311 In the precipitation frequency distribution, many models show a bimodal distribution in the
312 heavy precipitating tropical ocean region (Fig. 6b) but not in the light precipitating
313 subtropical ocean region (Fig. 6e) or the heavy precipitating mid-latitude land region (Fig.
314 6h). The bimodal frequency distribution is a commonly found in models and is seemingly
315 independent of resolution (e.g., Lin et al. 2013; Kooperman et al. 2018; Chen et al. 2021;
316 Ma et al. 2022; Martinez-Villalobos et al. 2022). It is not generally found in satellite-based
317 observational datasets, but this could be because the range of sensitivity to precipitation
318 rates is too narrow. Ma et al. (2022) compared the frequency distributions in AMIP and
319 HighResMIP (High Resolution Model Intercomparison Project, Haarsma et al. 2016) from



320 CMIP6 and DYAMOND (DYnamics of the Atmospheric general circulation Modeled On
321 Non-hydrostatic Domains, Satoh et al. 2019; Stevens et al. 2019) models, where they
322 showed that the bimodal frequency distribution appears in many AMIP (~100km),
323 HighResMIP (~50km), and even DYAMOND (~4km) models. Convective
324 parameterizations have been speculated as a cause of the light rain frequency peak (Lin
325 et al. 2013; Kooperman et al. 2018; Chen et al. 2021), though some models show a
326 convective precipitation peak at heavier precipitation than the peak from large-scale
327 precipitation (Martinez-Villalobos et al. 2022). ERA5 reanalysis also shows a bimodal
328 frequency distribution (Fig. 6b), which is not surprising considering that the reproduced
329 precipitation in ERA5 heavily depends on the model, thus exhibits this common model
330 behavior. Because of the heavy reliance on model physics to generate its precipitation
331 (as opposed to fields like wind, for which observations are directly assimilated), in this
332 study we do not include ERA5 precipitation among the observational products used for
333 model evaluation.

334

335 Based on the precipitation amount, frequency, and cumulative amount fraction curves,
336 we calculate 10 metrics (Amount peak, Amount P10, Amount P90, Frequency peak,
337 Frequency P10, Frequency P90, Unevenness, FracPRdays, SDII, and Perkins score) as
338 described in Section 3. Figure 7 shows the metrics with the modified IPCC AR6 climate
339 reference regions for satellite-based observations (black), ERA5 (gray), CMIP5 (blue),
340 and CMIP6 (red) models. The metric values vary widely across regions, especially in
341 Amount peak, Frequency peak, Unevenness, FracPRdays, and SDII, demonstrating the
342 additional detail provided by regional-scale precipitation-distribution metrics. In terms of



343 the metrics based on the amount distribution (Fig. 7a-c), many models tend to simulate
344 an Amount peak that is too light, an Amount P10 that is too high, and an Amount P90 that
345 is too low compared to the observations, moreso in oceanic regions (regions 47-62) than
346 in land regions. Similarly for the metrics based on the frequency distribution (Fig. 7d-f),
347 many models show light Frequency peaks, overestimated Frequency P10, and
348 underestimated Frequency P90 compared to observations. The similarity between
349 frequency distribution curves (i.e., Perkins score) is higher in land regions than in ocean
350 regions. Also, many models overestimate Unevenness and FracPRdays and
351 underestimate SDII. These results indicate that overall, models simulate more frequent
352 weak precipitation and less heavy precipitation compared to the observations, consistent
353 with many previous studies (e.g., Dai 2006; Pendergrass and Hartmann 2014; Trenberth
354 et al. 2017; Chen et al. 2021; Ma et al. 2022).

355

356 As expected from previous work, observations disagree substantially in some regions
357 (e.g., polar and high latitude regions) and/or for some metrics (e.g., Amount P90,
358 Frequency P90). In some cases the observational spread is much larger than that of the
359 models. We examine the observational discrepancy or spread by the ratio between the
360 standard deviation of the five satellite-based observations (IMERG, TRMM, CMORPH,
361 GPCP, PERSIANN) and the standard deviation of all CMIP 5 and 6 models (Fig. 8). The
362 standard deviation of observations is much larger near polar regions and high latitude
363 regions compared to the models' standard deviation for most metrics, as expected from
364 the orbital configurations of the most relevant satellite constellations for precipitation
365 (which exclude high latitudes). The Amount P90 and Frequency P90 metrics show the



366 largest observational discrepancy among the metrics, with standard deviations of 1.5 to
367 3 times larger over some high latitude regions and about 3-8 times larger over polar
368 regions in observations compared to the models. On the other hand, Unevenness,
369 FracPRdays, and Amount P10 show the least observational discrepancy – the models'
370 standard deviation is about 2-8 times larger than for observations over some tropical and
371 subtropical regions; nonetheless, the standard deviation among observations is larger
372 over most of the high latitude and polar regions. Model evaluation in the regions with large
373 disagreement among observational products remains a challenge. Note that the standard
374 deviation of five observations would be sensitive as there are outlier observations for
375 some regions and metrics (e.g., many ocean regions in Amount P90). Moreover,
376 observational uncertainties are rarely well quantified or understood, so agreements
377 among observational datasets may not always allow us to rule out common errors among
378 observations (e.g., for warm light precipitation over the subtropical ocean).

379

380 To attempt to account for discrepancies among observational datasets in the model
381 evaluation framework, we use two different approaches to evaluate model performance
382 with multiple observations, as described in Section 3.4. The first approach is to assess
383 the number of models that are within the observational range. Figure 9 shows the CMIP6
384 model evaluation with each metric, and the regions where the standard deviation among
385 observations is larger than among models are stippled gray to avoid them from the model
386 performance evaluation. In Amount peak, some subtropical regions (e.g., ARP, EAS,
387 NEPO, CAU, and WSAF) show relatively good model performance (more than 70% of
388 models fall in the observational range), while some tropical and subtropical (e.g., PITCZ,



389 AITCZ, and SEPO) and polar (e.g., RAR, EAN, and WAN) regions show poor model
390 performance (less than 30% of models fall in observational range). For Amount P10,
391 many regions are poorly captured by the simulations, except for some subtropical land
392 regions (e.g., EAS, NCA, CAU, and WSAF). In Amount P90, most regions are uncertain
393 (i.e., the standard deviation among observations is larger than among models) making it
394 difficult to evaluate model performance, while some tropical regions near the Indo-Pacific
395 warmpool (EIO, SEA, NWPO, and NAU) exhibit very good model performance (more than
396 90% of models fall in observational range). In the Frequency metrics (peak, P10, and
397 P90), more regions are difficult to evaluate model performance than in Amount metrics,
398 while in some tropical and subtropical regions (e.g., PITCZ, SWPO, NWPO, SEA, SAO,
399 and NES) model performance is good. However, good model performance could
400 alternatively arise from a large observational range (see Fig. 7). Unevenness,
401 FracPRdays, SDII, and Perkins score have a smaller fraction of models within the
402 observational range in tropical regions than the Amount and Frequency metrics. In
403 particular, fewer than 10% of CMIP6 models fall within the observational range for
404 Unevenness and FracPRdays over some tropical oceanic regions (e.g., PITCZ, NEPO,
405 SEPO, AITCZ, NEAO, SAO, and SIO).

406

407 Examining the fraction of CMIP5 models falling within the range of observations, CMIP5
408 models have a spatial pattern of model performance similar to that of CMIP6 models (see
409 Fig. S2 in supplement), and the improvement from CMIP5 to CMIP6 seems subtle. We
410 quantitatively assess the improvement from CMIP5 to CMIP6 by subtracting the
411 percentage of CMIP5 from CMIP6 models falling within the range of observations (Fig.



412 10). For some metrics (e.g., Amount peak, Amount and Frequency P10, and Amount and
413 Frequency P90) and for some tropical and subtropical regions (e.g., SEA, EAS, SAS,
414 ARP, and SAH), improvement is apparent. Compared to CMIP5, 5-25% more CMIP6
415 models fall in the observational range in these regions. However, for the other metrics
416 (e.g., Frequency peak, FracPRdays, SDII, Perkins score), CMIP6 models perform
417 somewhat worse. Over some tropical and subtropical oceanic regions (e.g., PITCZ,
418 NEPO, AITCZ, and NEAO), 5-25% more CMIP6 than CMIP5 models are out of the
419 observational range. This result is from all available CMIP5 and CMIP6 models, so it may
420 reflect the fact that some models are participated in only CMIP5 or CMIP6, but not both
421 (see Table 2). To isolate improvements that may have occurred between successive
422 generations of the same models, we also compared only the models that participated in
423 both CMIP5 and CMIP6 (see Fig. S3). Overall, the spatial characteristics of the
424 improvement/degradation in CMIP6 from CMIP5 is consistent, while more degradation is
425 apparent when we compare this subset of models, especially over the tropical oceanic
426 regions (e.g., PITCZ, AITCZ, NWPO, and SEPO).

427

428 The second approach to account for discrepancies among observations in model
429 performance evaluation is to count the number of models that are lower or higher than all
430 satellite-based observations for each metric and each region. Figure 11 shows the spatial
431 patterns of the model performance evaluation with each metric for CMIP6 models.
432 Underestimation is indicated by a negative sign, while overestimation is indicated by a
433 positive sign via the formulation described in Section 3.4. Amount peak is overall
434 underestimated in most regions, indicating the amount distributions in most CMIP6



435 models are shifted to lighter precipitation compared to observations. In many regions,
436 more than 50% of the CMIP6 models underestimate Amount peak. In particular, over
437 many tropical and southern hemisphere ocean regions (e.g., PITCZ, AITCZ, EIO, SEPO,
438 SAO, and SOO), more than 70% of the models underestimate the Amount peak. The
439 underestimation of Amount peak is accompanied by overestimation of Amount P10 and
440 underestimation of Amount P90. More than 70% of CMIP6 models overestimate Amount
441 P10 in many oceanic regions; especially in the southern and northern Pacific and Atlantic,
442 the southern Indian Ocean, and Southern Ocean more than 90% of the models
443 overestimate the observed Amount P10. For Amount P90, it appears that many models
444 fall within the observational range; however, observational range in Amount P90 (green
445 boxes in Fig. 7c) is large and driven primarily by just one observational dataset (IMERG),
446 especially in ocean regions.

447

448 For the frequency-based metrics (i.e., peak, P10, and P90; Figs. 11d-f), CMIP6 models
449 show similar bias characteristics to Amount metrics (Figs. 11a-c), although performance
450 is better than for Amount metrics. Over some tropical (e.g., NWPO, PITCZ, and SWPO)
451 and Eurasia (e.g., EEU, WSB, and ESB) regions, less than 10% of models fall outside of
452 the observed range. Unevenness and FracPRdays are severely overestimated in models.
453 More than 90% of models overestimate the observed Unevenness (Fig. 11g) and
454 FracPRdays (Fig. 11h) globally, especially over oceanic regions, consistent with
455 Pendergrass and Knutti (2018). Unevenness (i.e., number of the wettest days for the half
456 of annual precipitation) and FracPRdays (i.e., fraction of the number of annual
457 precipitating days above 1mm/day) are highly correlated to each other; correlations



458 between metrics will be discussed later. SDII is underestimated in many regions globally,
459 especially in some heavily-precipitating regions (e.g., PITCZ, AITCZ, EIO, SEA, NPO,
460 NAO, SWPO, and SOO). For the Perkins score, model simulations have poorer
461 performance in the tropics than in the mid-latitudes and polar regions. Performance by
462 these various metrics is generally consistent with the often-blamed too-frequent light
463 precipitation and too rare heavy precipitation in simulations.

464

465 The characteristics of CMIP5 compared to CMIP6 simulations (Fig. S4) show little
466 indication of improvement. Here we quantitatively evaluate the improvement in CMIP6
467 from CMIP5 for each metric and region. Figure 12 shows the difference between CMIP5
468 and CMIP6 in terms of the percentage of models that under- or over-estimate each metric.
469 In mid-latitudes, there appears to have been an improvement in performance, however in
470 the tropics, there appears to be more degradation. Over some heavily-precipitating
471 tropical regions (e.g., PITCZ, AITCZ, EIO, and NWPO), 10-25% more models in CMIP6
472 than in CMIP5 overestimate Amount P10, Unevenness, and FracPRdays and
473 underestimate/underperform on Amount peak, SDII, and Perkins score. This indicates
474 that CMIP6 models simulate more frequent light precipitation and less frequent heavy
475 precipitation over the heavily-precipitating tropical regions. Over some mid-latitude land
476 regions (e.g., EAS, ESB, RFE, and ENA), on the other hand, 5-20% more models in
477 CMIP6 than in CMIP5 simulate precipitation distributions close to observations (i.e., less
478 light precipitation and more heavy precipitation). To evaluate the improvement between
479 model generation, we also compare only the models that participated in both CMIP5 and
480 CMIP6 (Fig. S5) rather than all available CMIP5 and CMIP6 models. For the subset of



481 models participating in both generations, the improvement characteristics are similar for
482 all models, although more degradation is exhibited over some tropical oceanic regions
483 (e.g., PITCZ, NWPO, and SWPO). This also indicates that some models newly
484 participating in CMIP6, and not in the CMIP5, have higher than average performance.

485

486 4.3. Correlation between metrics

487 Each precipitation distribution metric implemented in this study is chosen to target
488 different aspects of the distribution of precipitation. To the extent that precipitation
489 probability distributions are governed by a small number of key parameters (as argued by
490 Martinez-Villalobos and Neelin 2019), we should expect additional metrics to be highly
491 correlated. Figure 13 shows the global weighted average of correlation coefficients
492 between the precipitation distribution metrics across CMIP5 and CMIP6 models. Higher
493 correlation coefficients are found to be between Amount P90 and Frequency P90 (0.98)
494 and between Amount P10 and Frequency P10 (0.67). This is expected because the
495 amount and frequency distributions differ only by a factor of the precipitation rate (e.g.,
496 Pendergrass and Hartmann 2014). Another higher correlation coefficient is between
497 Unevenness and FracPRdays (0.77), indicating that the number of the heaviest
498 precipitating days for half of annual precipitation and the total number of annual
499 precipitating days are related. Amount and Frequency peak metrics are negatively
500 correlated to P10 metrics and positively correlated to P90 metrics, but the correlation
501 coefficients are not very high (lower than 0.62). This is because the peak metrics focus
502 on typical precipitation, rather than the light and heavy ends of the distribution that are
503 the focus of P10 and P90 metrics. SDII is more negatively correlated with Amount P10 (-



504 0.67) and positively correlated with Amount peak (0.61) and less so with Amount P90
505 (0.48), implying that SDII is mainly influenced by weak precipitation amounts rather than
506 heavy precipitation amounts. The Perkins score shows relatively high negative correlation
507 with Unevenness (-0.62), FracPRdays (-0.59), and Amount P10 (-0.59). This indicates
508 that the discrepancy between the observed and modeled frequency distributions is partly
509 associated with the overestimated light precipitation in models. The correlation
510 coefficients between the metrics other than those discussed above are lower than 0.6.
511 While there is some redundant information within the collection of metrics included in our
512 framework, we retain all metrics so that others can select an appropriate subset for their
513 own application. This also preserves the ability to readily identify outlier behavior of an
514 individual model across a wide range of regions and statistics.

515

516 4.4. Influence of spatial resolution on metrics

517 Many metrics for the precipitation distribution are sensitive to the spatial resolution of
518 the underlying data (e.g., Pendergrass and Knutti 2018; Chen and Dai 2019). Figure 14
519 shows how our results (which are all based on data at 2° resolution) are impacted if we
520 calculate the metrics from data coarsened to 4° grid instead. As expected, there is clearly
521 some sensitivity to the spatial scale at which our precipitation distribution metrics are
522 computed, but the correlation among datasets (both models and observations) between
523 the two resolutions is very high, indicating that evaluations at either resolution should be
524 consistent. At the coarser resolution, Amount peak and SDII are consistently smaller (as
525 expected); Amount P10 and Frequency P10 tend to be smaller as well. Meanwhile,
526 Unevenness and FracPRdays are consistently large (as expected); Amount P90,



527 Frequency P90, and Perkins score are generally larger as well. Chen and Dai (2019)
528 discussed a grid aggregation effect that is associated with the increased probability of
529 precipitation as the horizontal resolution becomes coarser. This effect is clearly evident
530 with increased Unevenness (Fig. 14g), FracPRdays (Fig. 14h), and decreased SDII (Fig.
531 14i) in coarser resolution. However, despite these differences, the relative model
532 performance is not very sensitive to the spatial scale at which we apply our analysis. The
533 correlation coefficients between results based on all data interpolated to 2° or 4°
534 horizontal resolutions are above 0.9 for all of our distribution metrics. Conclusions on
535 model performance are relatively insensitive to the target resolution.

536

537

538 **5. Discussion**

539 Analyzing the distribution of precipitation intensity lags behind temperature and even
540 mean precipitation. Challenges include choosing appropriate metrics and analysis
541 resolution to characterize this highly non-gaussian variable and interpreting model skills
542 in the face of substantial observational uncertainty. Comparing results derived at 2° and
543 4° horizontal resolution for CMIP class models, we find that the quantitative changes in
544 assessed performance are highly consistent across models and consequently have little
545 impact on our conclusions. More work is needed to determine how suitable this collection
546 of metrics may be for evaluating models with substantially higher resolutions (e.g.,
547 HighResMIP, Haarsma et al. 2016). We note that more complex measures have been
548 designed to be scale independent (e.g., Martinez-Villalobos and Neelin 2019; Martinez-



549 Villalobos et al. 2022), and these may become increasingly important with continued
550 interest in models developed at substantially higher resolution.

551

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

560

561 The common model biases identified in this study are mainly associated with the
562 overestimated light precipitation and underestimated heavy precipitation. These biases
563 persist from deficiencies identified in earlier generation models (e.g., Dai 2006), and as
564 shown in this study there has been little improvement. One reason may be that these key
565 characteristics of precipitation are not commonly considered in the model development
566 process. Enabling modelers to more readily objectively evaluate simulated precipitation
567 distributions could perhaps serve as a guide to improvement. The current study aims to
568 provide a framework for objective evaluation of simulated precipitation distributions at
569 regional scales.

570



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

586

587

588 **6. Conclusion**

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

593



594 To facilitate the regional scale for evaluation, regions where precipitation characteristics
595 are relatively homogenous are identified. Our reference regions consist of existing IPCC
596 AR6 climate reference regions, with additional subdivisions based on homogeneity
597 analysis performed on precipitation distributions within each region. We partition the
598 global domain into heavy, moderate, and light precipitation regions using K-means
599 clustering of IMERG precipitation frequency and amount distributions. Our clustering
600 analysis reveals that the IPCC AR6 land regions are reasonably homogeneous in
601 precipitation character, while some ocean regions are relatively inhomogeneous,
602 including large portions of both heavy and light precipitating areas. To define more
603 homogeneous regions for the analysis of precipitation distributions, we have modified
604 some ocean regions to better fit the clustering results while retaining the original IPCC
605 AR6 land regions. The homogeneity between the region-averaged distribution and each
606 grid cell's distribution over the region is assessed by the three distinct similarity metrics
607 (Perkins score, K-S test, and A-D test). The homogeneity is overall improved in the
608 modified IPCC AR6 ocean regions. Although the clustering regions are obtained based
609 on the IMERG annual precipitation, the improved homogeneity is fairly consistent across
610 different datasets (TRMM, CMORPH, GPCP, PERSIANN, and ERA5) and seasons (MAM,
611 JJA, SON, and DJF). Use of these more homogeneous regions enables us to extract
612 more robust quantitative information from the distributions in each region.

613

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



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

629

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



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

650

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

660

661



662 **Code Availability**

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

671

672

673 **Data Availability**

674 All of the data used in this study are publicly available. The satellite-based precipitation
675 products used in this study (IMERG, TRMM, CMORPH, GPCP, and PERSIANN) and
676 ERA5 precipitation product are available on the Obs4MIPs at [https://esgf-](https://esgf-node.llnl.gov/projects/obs4mips/)
677 [node.llnl.gov/projects/obs4mips/](https://esgf-node.llnl.gov/projects/obs4mips/). The CMIP data is available on the ESGF at [node.llnl.gov/projects/esgf-llnl](https://esgf-
678 <a href=).

679

680

681 **Author contribution**

682 PG and AP designed the initial idea of the precipitation benchmarking framework. MA,
683 PU, PG, and JL advanced the idea and developed the framework. MA performed analysis.



684 MA, JL, and AO implemented the framework code into the PCMDI metrics package. MA
685 prepared the manuscript with contributions from all co-authors.

686

687

688 **Competing interests**

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

690

691

692 **Disclaimer**

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

705

706



707 **Acknowledgements**

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

722

723



724 **References**

- 725 Ahn, M., and I. Kang, 2018: A practical approach to scale-adaptive deep convection
726 in a GCM by controlling the cumulus base mass flux. *npj Clim. Atmos. Sci.*, **1**,
727 13, <https://doi.org/10.1038/s41612-018-0021-0>.
- 728 Ahn, M.-S., P. J. Gleckler, J. Lee, A. G. Pendergrass, and C. Jakob, 2022:
729 Benchmarking Simulated Precipitation Variability Amplitude across Time Scales.
730 *J. Clim.*, **35**, 3173–3196, <https://doi.org/10.1175/JCLI-D-21-0542.1>.
- 731 Ashouri, H., K. L. Hsu, S. Sorooshian, D. K. Braithwaite, K. R. Knapp, L. D. Cecil, B.
732 R. Nelson, and O. P. Prat, 2015: PERSIANN-CDR: Daily precipitation climate
733 data record from multisatellite observations for hydrological and climate studies.
734 *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)
735 00068.1.
- 736 Chakravarti, I. M., R. G. Laha, and J. Roy, 1967: Handbook of Methods of Applied
737 Statistics, Volume I: Techniques of Computation, Descriptive Methods, and
738 Statistical Inference. *John Wiley Sons*, 392–394.
- 739 Chen, D., and A. Dai, 2019: Precipitation Characteristics in the Community
740 Atmosphere Model and Their Dependence on Model Physics and Resolution. *J.*
741 *Adv. Model. Earth Syst.*, **11**, 2352–2374,
742 <https://doi.org/10.1029/2018MS001536>.



- 743 Chen, D., A. Dai, and A. Hall, 2021: The Convective-To-Total Precipitation Ratio and
744 the “Drizzling” Bias in Climate Models. *J. Geophys. Res. Atmos.*, **126**, 1–17,
745 <https://doi.org/10.1029/2020JD034198>.
- 746 Dai, A., 2006: Precipitation characteristics in eighteen coupled climate models. *J.*
747 *Clim.*, **19**, 4605–4630, <https://doi.org/10.1175/JCLI3884.1>.
- 748 Eyring, V., S. Bony, G. A. Meehl, C. A. Senior, B. Stevens, R. J. Stouffer, and K. E.
749 Taylor, 2016: Overview of the Coupled Model Intercomparison Project Phase 6
750 (CMIP6) experimental design and organization. *Geosci. Model Dev.*, **9**, 1937–
751 1958, <https://doi.org/10.5194/gmd-9-1937-2016>.
- 752 Fiedler, S., and Coauthors, 2020: Simulated Tropical Precipitation Assessed across
753 Three Major Phases of the Coupled Model Intercomparison Project (CMIP). *Mon.*
754 *Weather Rev.*, **148**, 3653–3680, <https://doi.org/10.1175/MWR-D-19-0404.1>.
- 755 Gleckler, P., C. Doutriaux, P. Durack, K. Taylor, Y. Zhang, D. Williams, E. Mason,
756 and J. Servonnat, 2016: A More Powerful Reality Test for Climate Models. *Eos*
757 (*Washington. DC.*), **97**, 20–24, <https://doi.org/10.1029/2016EO051663>.
- 758 Gleckler, P. J., K. E. Taylor, and C. Doutriaux, 2008: Performance metrics for climate
759 models. *J. Geophys. Res. Atmos.*, **113**, 1–20,
760 <https://doi.org/10.1029/2007JD008972>.
- 761 Hersbach, H., and Coauthors, 2020: The ERA5 global reanalysis. *Q. J. R. Meteorol.*
762 *Soc.*, **146**, 1999–2049, <https://doi.org/10.1002/qj.3803>.



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

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

769 Huffman, G. J., R. F. Adler, M. M. Morrissey, D. T. Bolvin, S. Curtis, R. Joyce, B.
770 McGavock, and J. Susskind, 2001: Global Precipitation at One-Degree Daily
771 Resolution from Multisatellite Observations. *J. Hydrometeorol.*, **2**, 36–50,
772 [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).

773 Iturbide, M., and Coauthors, 2020: An update of IPCC climate reference regions for
774 subcontinental analysis of climate model data: definition and aggregated
775 datasets. *Earth Syst. Sci. Data*, **12**, 2959–2970, [https://doi.org/10.5194/essd-12-
776 2959-2020](https://doi.org/10.5194/essd-12-2959-2020).

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

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



784 Kooperman, G. J., M. S. Pritchard, M. A. Burt, M. D. Branson, and D. A. Randall,
785 2016: Robust effects of cloud superparameterization on simulated daily rainfall
786 intensity statistics across multiple versions of the Community Earth System
787 Model. *J. Adv. Model. Earth Syst.*, **8**, 140–165,
788 <https://doi.org/10.1002/2015MS000574>.

789 Kooperman, G. J., M. S. Pritchard, T. A. O'Brien, and B. W. Timmermans, 2018:
790 Rainfall From Resolved Rather Than Parameterized Processes Better
791 Represents the Present-Day and Climate Change Response of Moderate Rates
792 in the Community Atmosphere Model. *J. Adv. Model. Earth Syst.*, **10**, 971–988,
793 <https://doi.org/10.1002/2017MS001188>.

794 Leung, L. R., and Coauthors, 2022: Exploratory Precipitation Metrics: Spatiotemporal
795 Characteristics, Process-Oriented, and Phenomena-Based. *J. Clim.*, **35**, 3659–
796 3686, <https://doi.org/10.1175/JCLI-D-21-0590.1>.

797 Lin, Y., M. Zhao, Y. Ming, J.-C. Golaz, L. J. Donner, S. A. Klein, V. Ramaswamy, and
798 S. Xie, 2013: Precipitation Partitioning, Tropical Clouds, and Intraseasonal
799 Variability in GFDL AM2. *J. Clim.*, **26**, 5453–5466, [https://doi.org/10.1175/JCLI-](https://doi.org/10.1175/JCLI-D-12-00442.1)
800 [D-12-00442.1](https://doi.org/10.1175/JCLI-D-12-00442.1).

801 Ma, Hsi-Yen, Stephen A. Klein, Jiwoo Lee, Min-Seop Ahn, Cheng Tao and Peter J.
802 Gleckler, 2022: Superior daily and sub-daily precipitation statistics for intense
803 and long-lived storms in global storm-resolving models. *Geophys. Res. Lett.*, in
804 revision.



- 805 MacQueen, J. B., 1967: Some methods for classification and analysis of multivariate
806 observations. *Berkeley Symp. Math. Stat. Probab.*, **VOL. 5.1**, 281–297.
- 807 Martinez-Villalobos, C., and J. D. Neelin, 2019: Why Do Precipitation Intensities Tend
808 to Follow Gamma Distributions? *J. Atmos. Sci.*, **76**, 3611–3631,
809 <https://doi.org/10.1175/JAS-D-18-0343.1>.
- 810 Martinez-Villalobos, C., J. D. Neelin, and A. G. Pendergrass, 2022: Metrics for
811 Evaluating CMIP6 Representation of Daily Precipitation Probability Distributions.
812 *J. Clim.*, 1–79, <https://doi.org/10.1175/JCLI-D-21-0617.1>.
- 813 Meehl, G. A., C. Covey, B. McAvaney, M. Latif, and R. J. Stouffer, 2005: Overview
814 of the Coupled Model Intercomparison Project. *Bull. Am. Meteorol. Soc.*, **86**, 89–
815 96, <https://doi.org/10.1175/BAMS-86-1-89>.
- 816 Meehl, G. A., C. Covey, T. Delworth, M. Latif, B. McAvaney, J. F. B. Mitchell, R. J.
817 Stouffer, and K. E. Taylor, 2007: THE WCRP CMIP3 Multimodel Dataset: A New
818 Era in Climate Change Research. *Bull. Am. Meteorol. Soc.*, **88**, 1383–1394,
819 <https://doi.org/10.1175/BAMS-88-9-1383>.
- 820 Meehl, G. A., G. J. Boer, C. Covey, M. Latif, and R. J. Stouffer, 2000: The Coupled
821 Model Intercomparison Project (CMIP). *Bull. Am. Meteorol. Soc.*, **81**, 313–318,
822 [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).
- 823 Pendergrass, A. G., and C. Deser, 2017: Climatological Characteristics of Typical
824 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)
825 [0684.1](https://doi.org/10.1175/JCLI-D-16-0684.1).



- 826 Pendergrass, A. G., and D. L. Hartmann, 2014: Two Modes of Change of the
827 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)
828 [14-00182.1](https://doi.org/10.1175/JCLI-D-14-00182.1).
- 829 Pendergrass, A. G., and R. Knutti, 2018: The Uneven Nature of Daily Precipitation
830 and Its Change. *Geophys. Res. Lett.*, **45**, 11,980-11,988,
831 <https://doi.org/10.1029/2018GL080298>.
- 832 Pendergrass, A. G., P. J. Gleckler, L. R. Leung, and C. Jakob, 2020: Benchmarking
833 Simulated Precipitation in Earth System Models. *Bull. Am. Meteorol. Soc.*, **101**,
834 E814–E816, <https://doi.org/10.1175/BAMS-D-19-0318.1>.
- 835 Perkins, S. E., A. J. Pitman, N. J. Holbrook, and J. McAneney, 2007: Evaluation of
836 the AR4 Climate Models' Simulated Daily Maximum Temperature, Minimum
837 Temperature, and Precipitation over Australia Using Probability Density
838 Functions. *J. Clim.*, **20**, 4356–4376, <https://doi.org/10.1175/JCLI4253.1>.
- 839 Roca, R., L. V. Alexander, G. Potter, M. Bador, R. Jucá, S. Contractor, M. G.
840 Bosilovich, and S. Cloché, 2019: FROGS: a daily 1° × 1° gridded precipitation
841 database of rain gauge, satellite and reanalysis products. *Earth Syst. Sci. Data*,
842 **11**, 1017–1035, <https://doi.org/10.5194/essd-11-1017-2019>.
- 843 Stephens, M. A., 1974: EDF Statistics for Goodness of Fit and Some Comparisons.
844 *J. Am. Stat. Assoc.*, **69**, 730–737, <https://doi.org/10.2307/2286009>.
- 845 Sun, Y., S. Solomon, A. Dai, and R. W. Portmann, 2006: How Often Does It Rain? *J.*
846 *Clim.*, **19**, 916–934, <https://doi.org/10.1175/JCLI3672.1>.



- 847 Sun, Y., S. Solomon, A. Dai, and R. W. Portmann, 2007: How Often Will It Rain? *J.*
848 *Clim.*, **20**, 4801–4818, <https://doi.org/10.1175/JCLI4263.1>.
- 849 Swenson, L. M., and R. Grotjahn, 2019: Using Self-Organizing Maps to Identify
850 Coherent CONUS Precipitation Regions. *J. Clim.*, **32**, 7747–7761,
851 <https://doi.org/10.1175/JCLI-D-19-0352.1>.
- 852 Tang, S., P. Gleckler, S. Xie, J. Lee, M.-S. Ahn, C. Covey, and C. Zhang, 2021:
853 Evaluating Diurnal and Semi-Diurnal Cycle of Precipitation in CMIP6 Models
854 Using Satellite- and Ground-Based Observations. *J. Clim.*, 1–56,
855 <https://doi.org/10.1175/JCLI-D-20-0639.1>.
- 856 Taylor, K. E., R. J. Stouffer, and G. A. Meehl, 2012: An overview of CMIP5 and the
857 experiment design. *Bull. Amer. Meteor. Soc.*, **93**, 485–498,
858 <https://doi.org/10.1175/BAMS-D-11-00094.1>.
- 859 Trenberth, K. E., A. Dai, R. M. Rasmussen, and D. B. Parsons, 2003: The Changing
860 Character of Precipitation. *Bull. Am. Meteorol. Soc.*, **84**, 1205–1218,
861 <https://doi.org/10.1175/BAMS-84-9-1205>.
- 862 Trenberth, K. E., and Y. Zhang, 2018: How Often Does It Really Rain? *Bull. Am.*
863 *Meteorol. Soc.*, **99**, 289–298, <https://doi.org/10.1175/BAMS-D-17-0107.1>.
- 864 Trenberth, K. E., Y. Zhang, and M. Gehne, 2017: Intermittency in Precipitation:
865 Duration, Frequency, Intensity, and Amounts Using Hourly Data. *J.*
866 *Hydrometeorol.*, **18**, 1393–1412, <https://doi.org/10.1175/JHM-D-16-0263.1>.



- 867 U.S. DOE. 2020. Benchmarking Simulated Precipitation in Earth System Models
868 Workshop Report, DOE/SC-0203, U.S. Department of Energy Office of Science,
869 Biological and Environmental Research (BER) Program. Germantown,
870 Maryland, USA.
- 871 Waliser, D., and Coauthors, 2020: Observations for Model Intercomparison Project
872 (Obs4MIPs): status for CMIP6. *Geosci. Model Dev.*, **13**, 2945–2958,
873 <https://doi.org/10.5194/gmd-13-2945-2020>.
- 874 Wei, G., H. Lü, W. T. Crow, Y. Zhu, J. Wang, and J. Su, 2017: Evaluation of Satellite-
875 Based Precipitation Products from IMERG V04A and V03D, CMORPH and
876 TMPA with Gauged Rainfall in Three Climatologic Zones in China. *Remote
877 Sens.*, **10**, 30, <https://doi.org/10.3390/rs10010030>.
- 878 Xie, P., R. Joyce, S. Wu, S. H. Yoo, Y. Yarosh, F. Sun, and R. Lin, 2017:
879 Reprocessed, bias-corrected CMORPH global high-resolution precipitation
880 estimates from 1998. *J. Hydrometeorol.*, **18**, 1617–1641,
881 <https://doi.org/10.1175/JHM-D-16-0168.1>.
- 882 Zhang, C., X. Chen, H. Shao, S. Chen, T. Liu, C. Chen, Q. Ding, and H. Du, 2018:
883 Evaluation and intercomparison of high-resolution satellite precipitation
884 estimates-GPM, TRMM, and CMORPH in the Tianshan Mountain Area. *Remote
885 Sens.*, **10**, <https://doi.org/10.3390/rs10101543>.
- 886 Zhang, X., L. Alexander, G. C. Hegerl, P. Jones, A. K. Tank, T. C. Peterson, B. Trewin,
887 and F. W. Zwiers, 2011: Indices for monitoring changes in extremes based on



888 daily temperature and precipitation data. *Wiley Interdiscip. Rev. Clim. Chang.*, **2**,

889 851–870, <https://doi.org/10.1002/wcc.147>.

890



891 **Tables**

892

893

894

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

896

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)

897

898

899

900

901

902

903

904

905

906

907

908

909



910

911 Table 2. CMIP5 and CMIP6 models used in this study and their horizontal resolution. The

912 number in parentheses indicates the number of realizations used for each model. Note

913 that the horizontal resolution information is obtained from the number of grids, and it may

914 vary slightly if the grid interval is not linear.

915

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

916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940



941

942 Table 3. Precipitation distribution metrics implemented in this study.

943

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

944

945

946

947

948

949



950
 951 Table 4. List of climate reference regions used in this study. The new ocean regions
 952 defined in this study are highlighted in bold.
 953

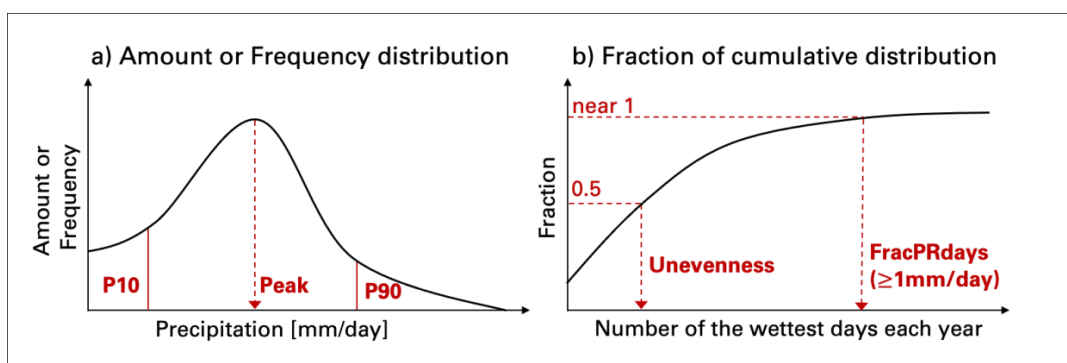
1	GIC	Greenland/Iceland	22	WAF	Western-Africa	43	SAU	S.Australia
2	NWN	N.W.North-America	23	CAF	Central-Africa	44	NZ	New-Zealand
3	NEN	N.E.North-America	24	NEAF	N.Eastern-Africa	45	EAN	E.Antarctica
4	WNA	W.North-America	25	SEAF	S.Eastern-Africa	46	WAN	W.Antarctica
5	CNA	C.North-America	26	WSAF	W.Southern-Africa	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	NPO	N.Pacific-Ocean
11	NSA	N.South-America	32	RFE	Russian-Far-East	53	NWP O	N.W.Pacific-Ocean
12	NES	N.E.South-America	33	WCA	W.C.Asia	54	NEPO	N.E.Pacific-Ocean
13	SAM	South-American-Monsoon	34	ECA	E.C.Asia	55	PITCZ	Pacific-ITCZ
14	SWS	S.W.South-America	35	TIB	Tibetan-Plateau	56	SWPO	S.W.Pacific-Ocean
15	SES	S.E.South-America	36	EAS	E.Asia	57	SEPO	S.E.Pacific-Ocean
16	SSA	S.South-America	37	ARP	Arabian-Peninsula	58	NAO	N.Atlantic-Ocean
17	NEU	N.Europe	38	SAS	S.Asia	59	NEAO	N.E.Atlantic-Ocean
18	WCE	West&Central-Europe	39	SEA	S.E.Asia	60	AITCZ	Atlantic-ITCZ
19	EEU	E.Europe	40	NAU	N.Australia	61	SAO	S.Atlantic-Ocean
20	MED	Mediterranean	41	CAU	C.Australia	62	SOO	Southern-Ocean
21	SAH	Sahara	42	EAU	E.Australia			

954



955 **Figures**

956
957
958
959
960
961

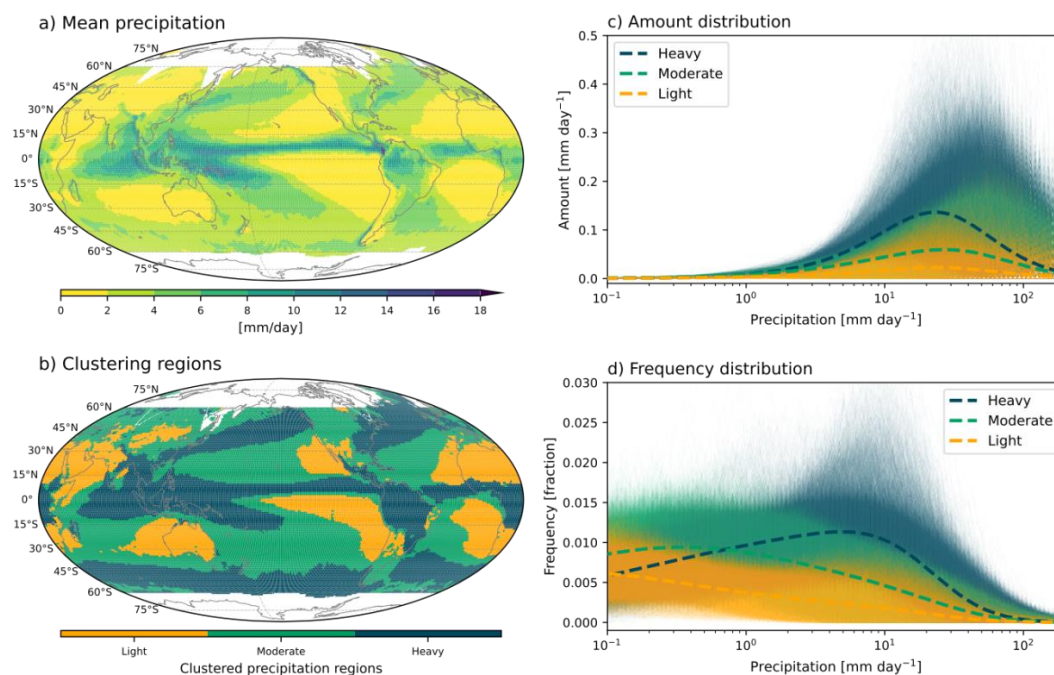


962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986

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



987
988
989

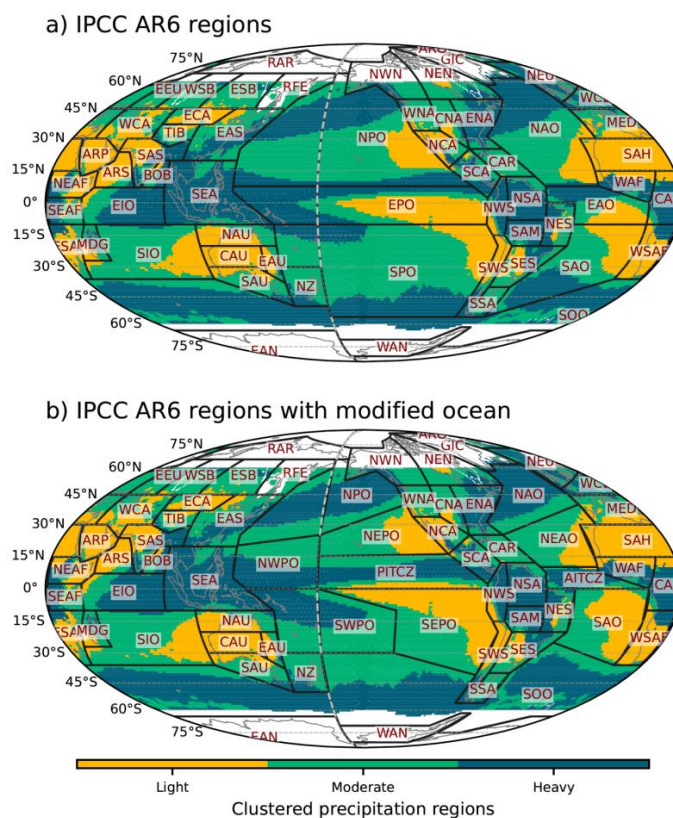


990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009

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.



1010
1011
1012

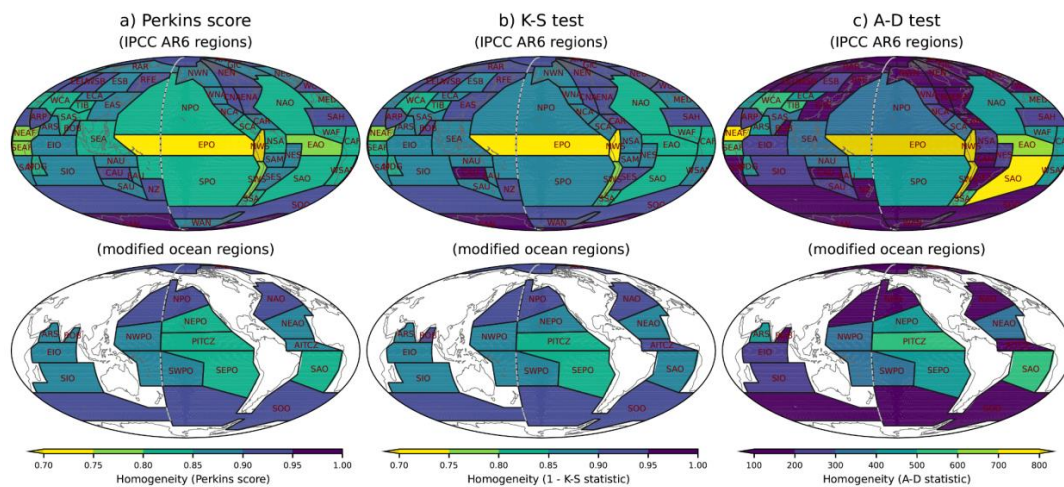


1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028

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.



1029
1030
1031

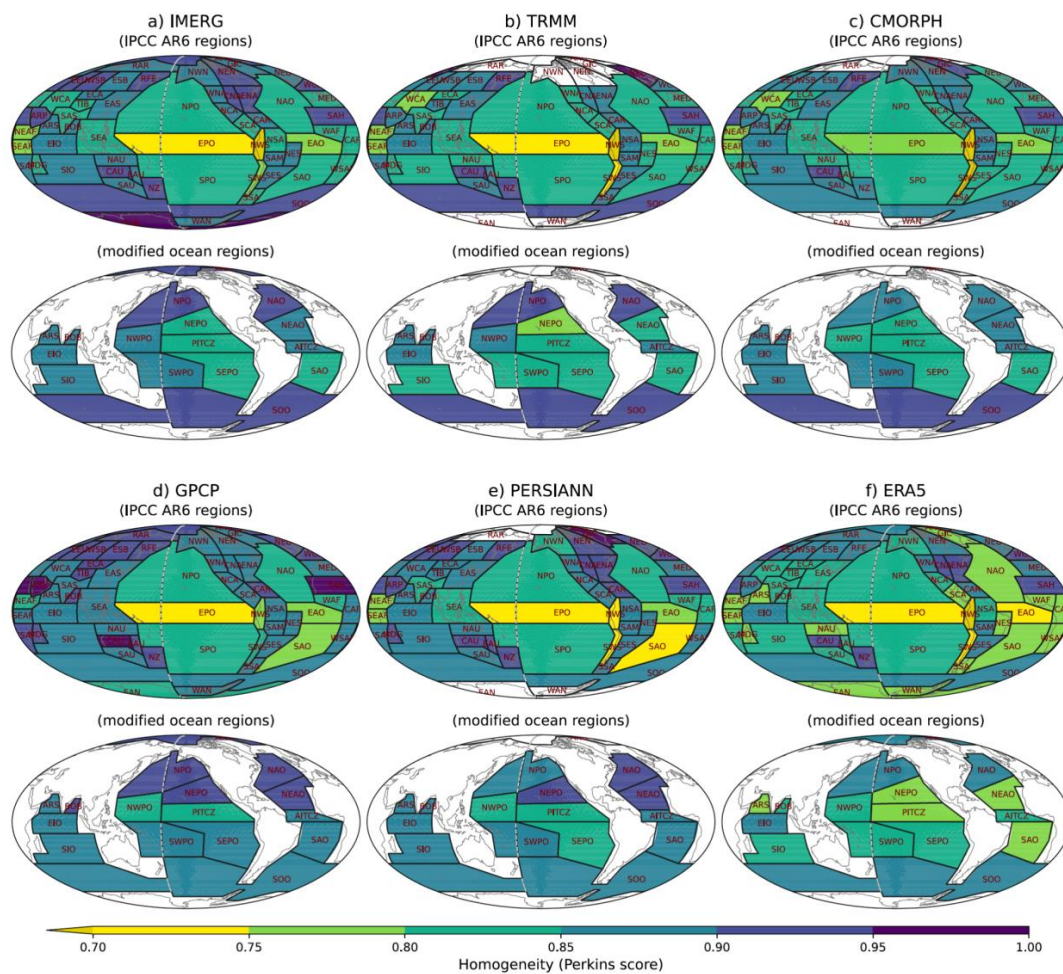


1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056

Figure 4. Homogeneity estimated by a) Perkins score, b) K-S test, and c) A-D test between the region averaged and each grid's frequency distributions of IMERG precipitation for the IPCC AR6 climate reference regions (upper) and the modified ocean regions (bottom). Darker color indicates higher homogeneity across all panels.



1057
1058
1059

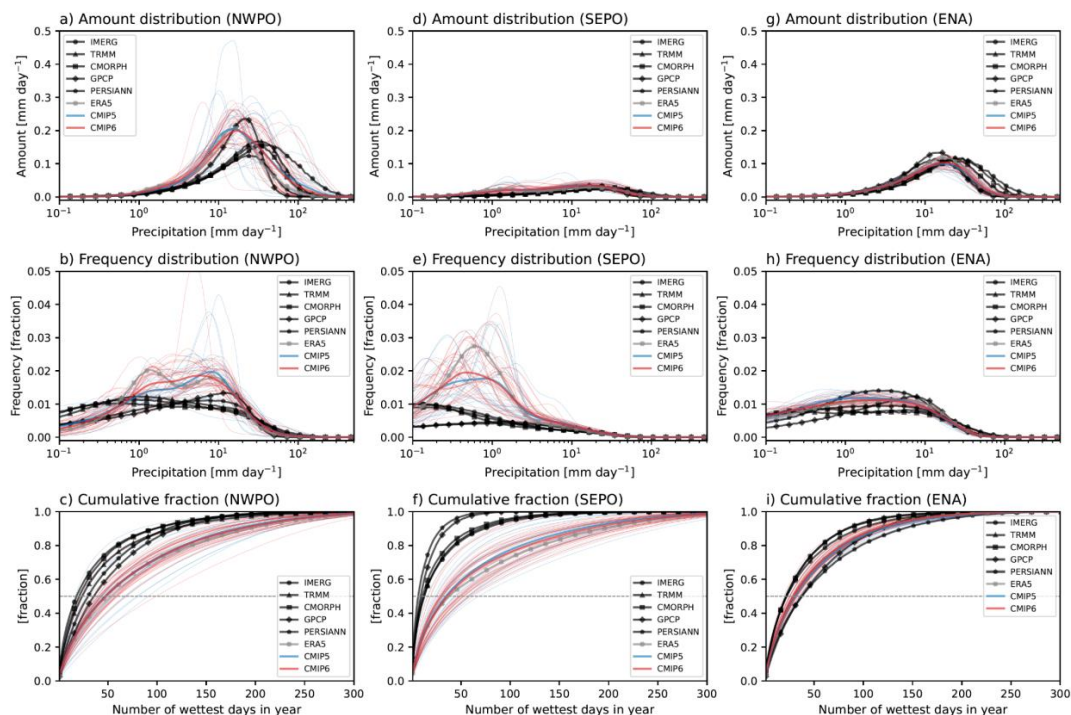


1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071

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



1072
 1073
 1074

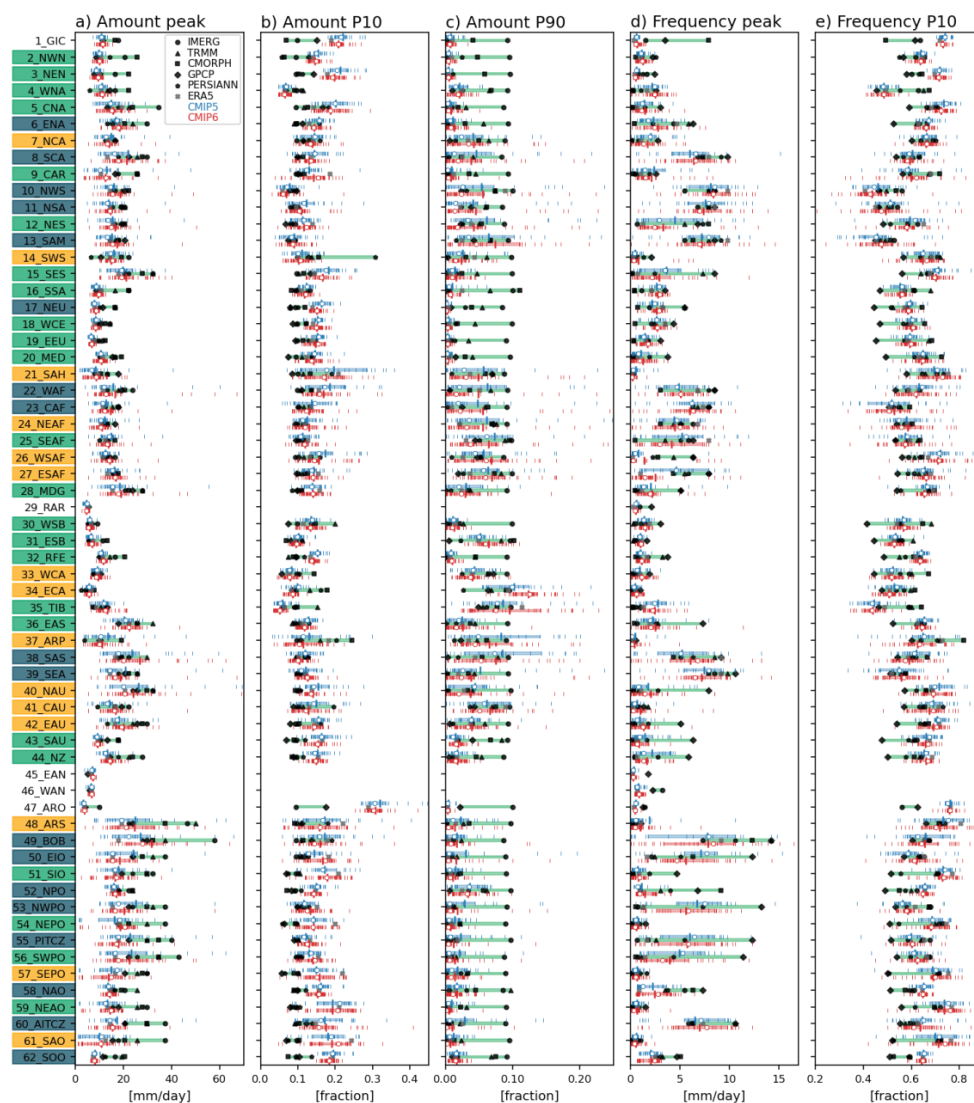


1075
 1076
 1077
 1078
 1079
 1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093

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.



1094
 1095

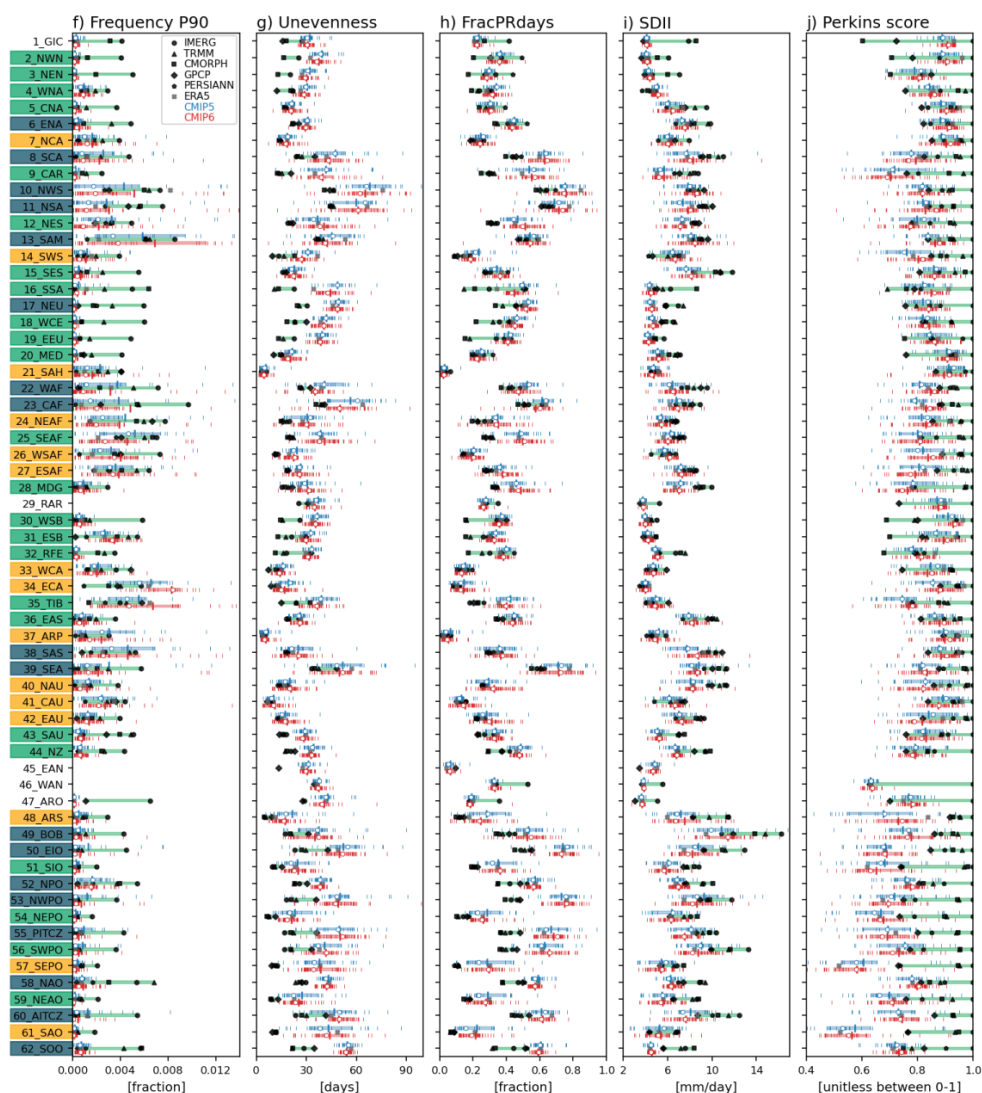


1096
 1097

1098 Figure 7. Precipitation distribution metrics for a) Amount peak, b) Amount P10, c)
 1099 Amount P90, d) Frequency peak, e) Frequency P10, f) Frequency P90, g) Unevenness,
 1100 h) FracPRdays, i) SDII, and j) Perkins score over the modified IPCC AR6 regions.
 1101 Black, gray, blue, and red curves indicate the satellite-based observations, reanalysis,
 1102 CMIP5 models, and CMIP6 modes, respectively. Thin and thick vertical marks for CMIP
 1103 models respectively indicate distributions for each model and multi-model average.
 1104 Open circle mark for CMIP models indicates the multi-model median. Green shade



1105 represents the range between the minimum and maximum values of satellite-based
 1106 observations. Blue and red shades respectively represent the range between 25th and
 1107 75th model values for CMIP 5 and 6 models. Y-axis labels are shaded with the three
 1108 colors as the same in Fig. 2b, indicating dominant precipitating characteristics. Note that
 1109 regions 1-46 are land and land-ocean mixed regions, and 47-62 are ocean regions.
 1110
 1111

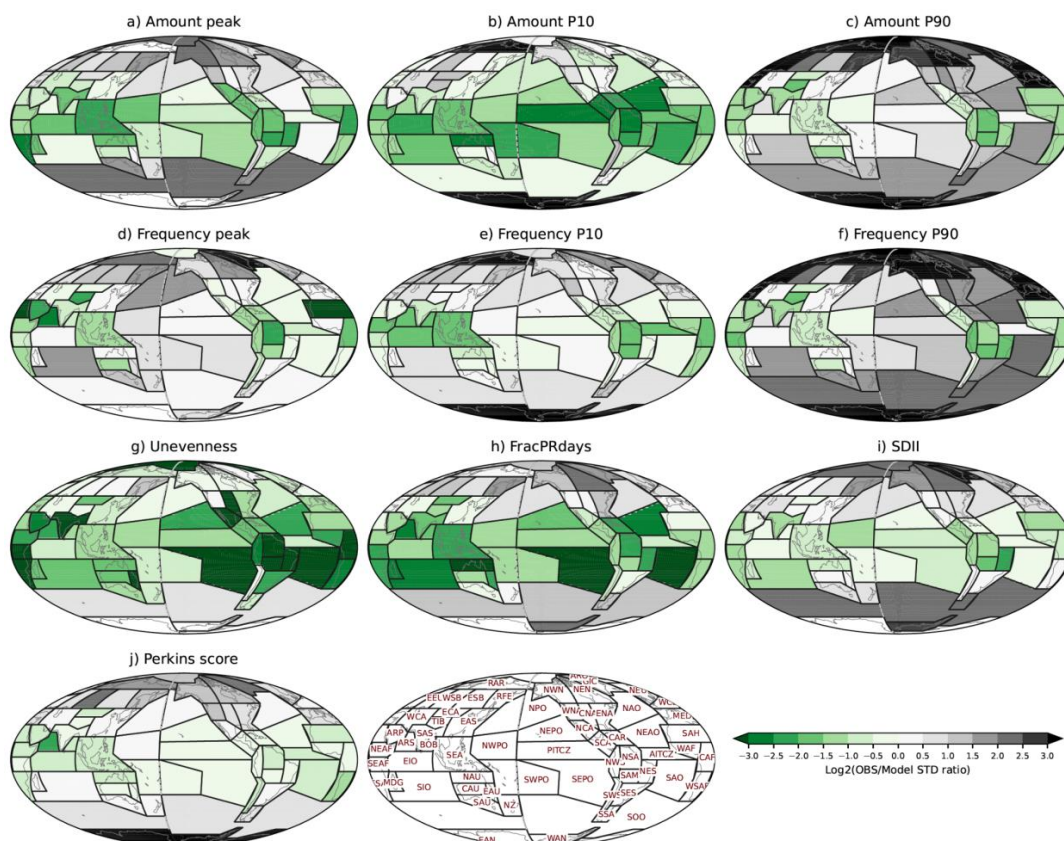


1112
 1113
 1114
 1115

Figure 7. (continued)



1116
1117

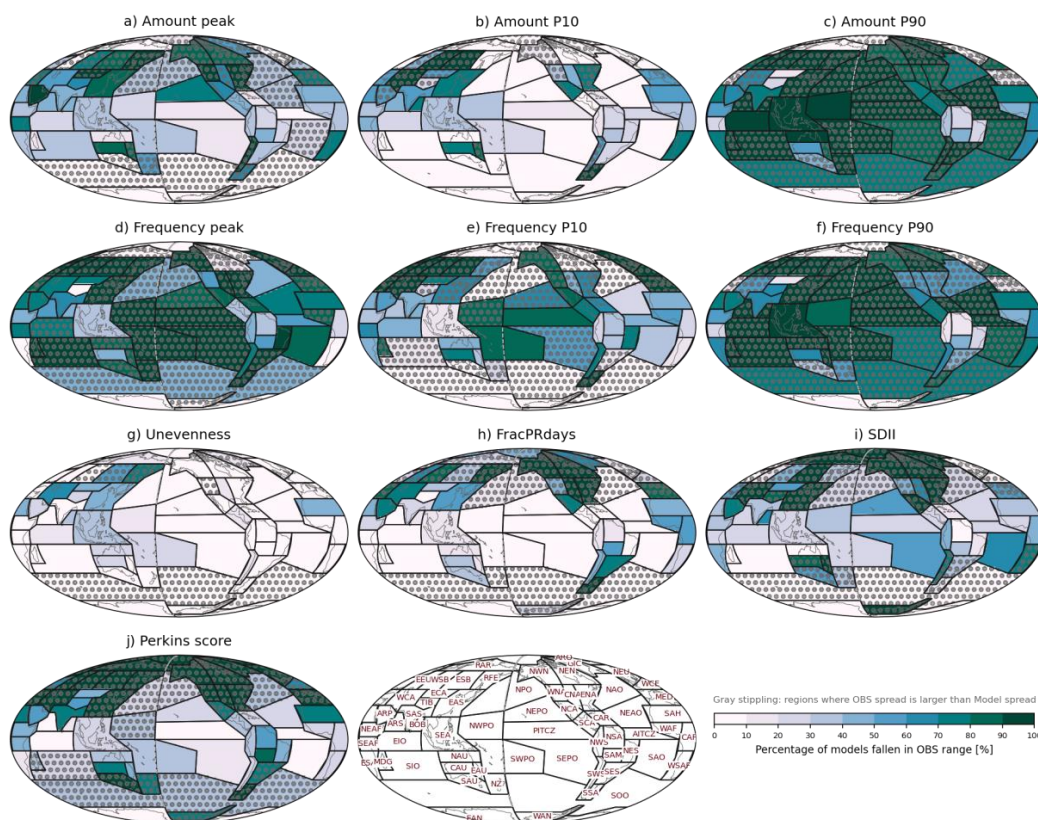


1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

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.



1134
1135

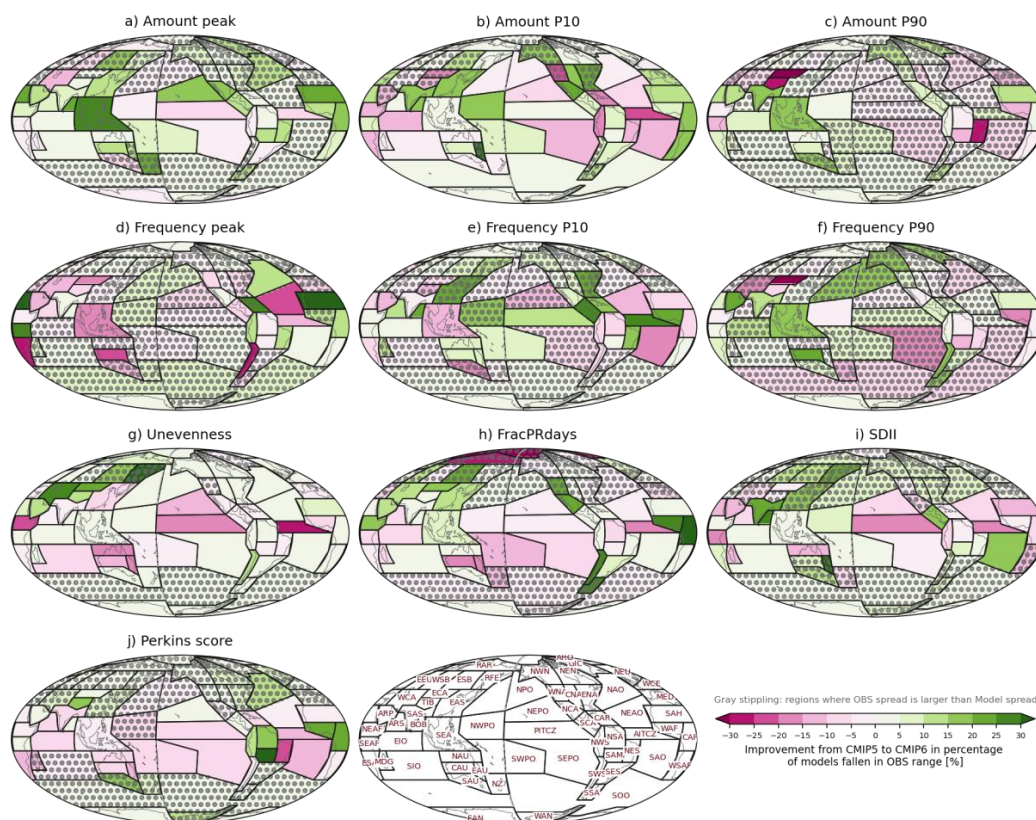


1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152

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.



1153
1154

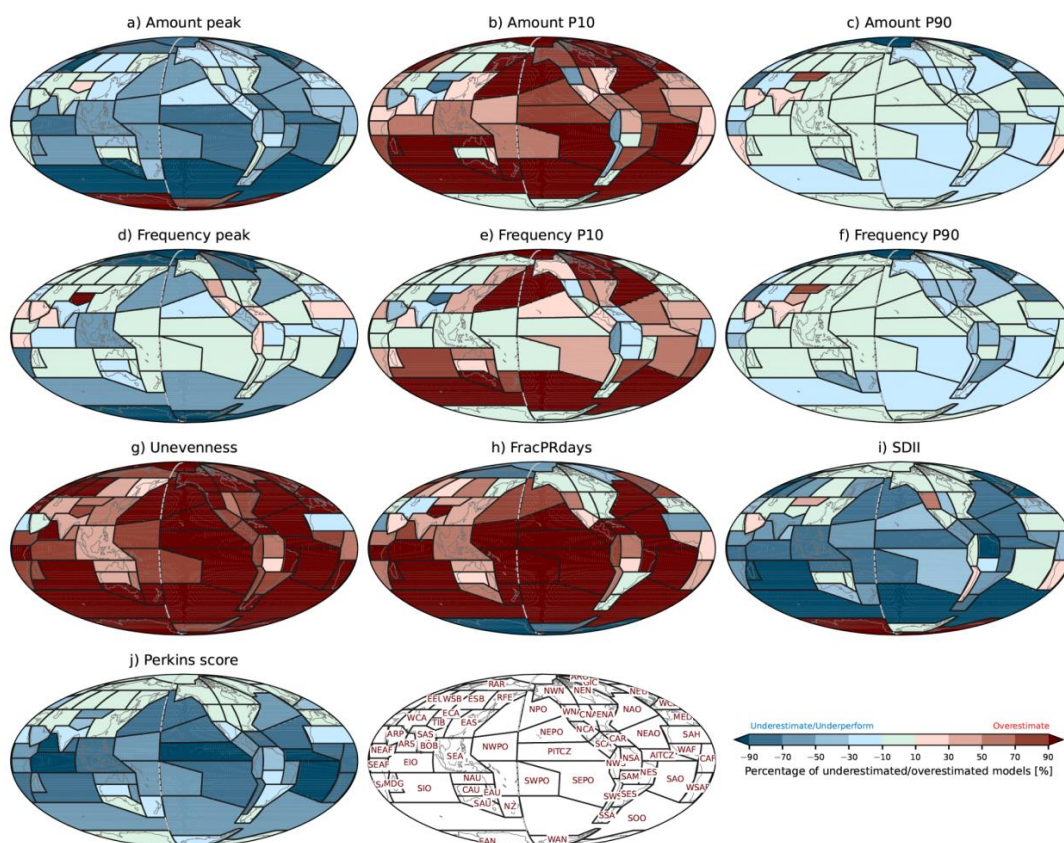


1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171

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.



1172
1173

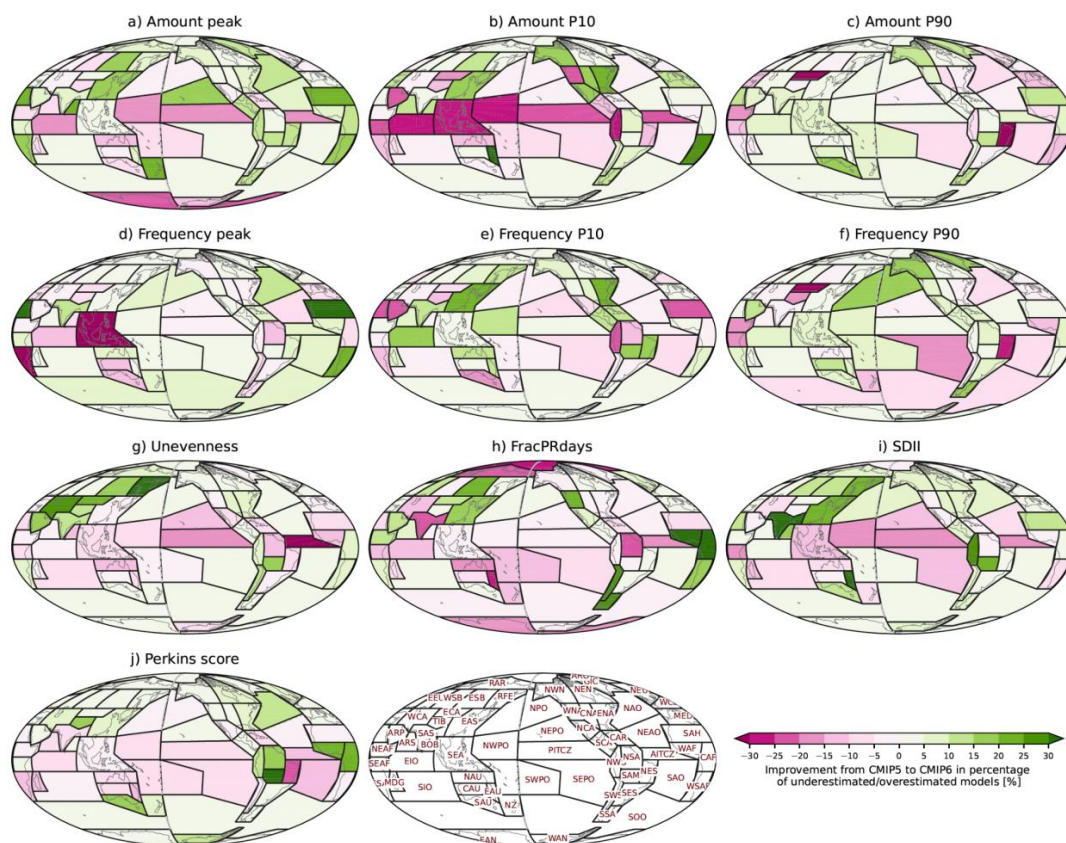


1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189

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.



1190
1191

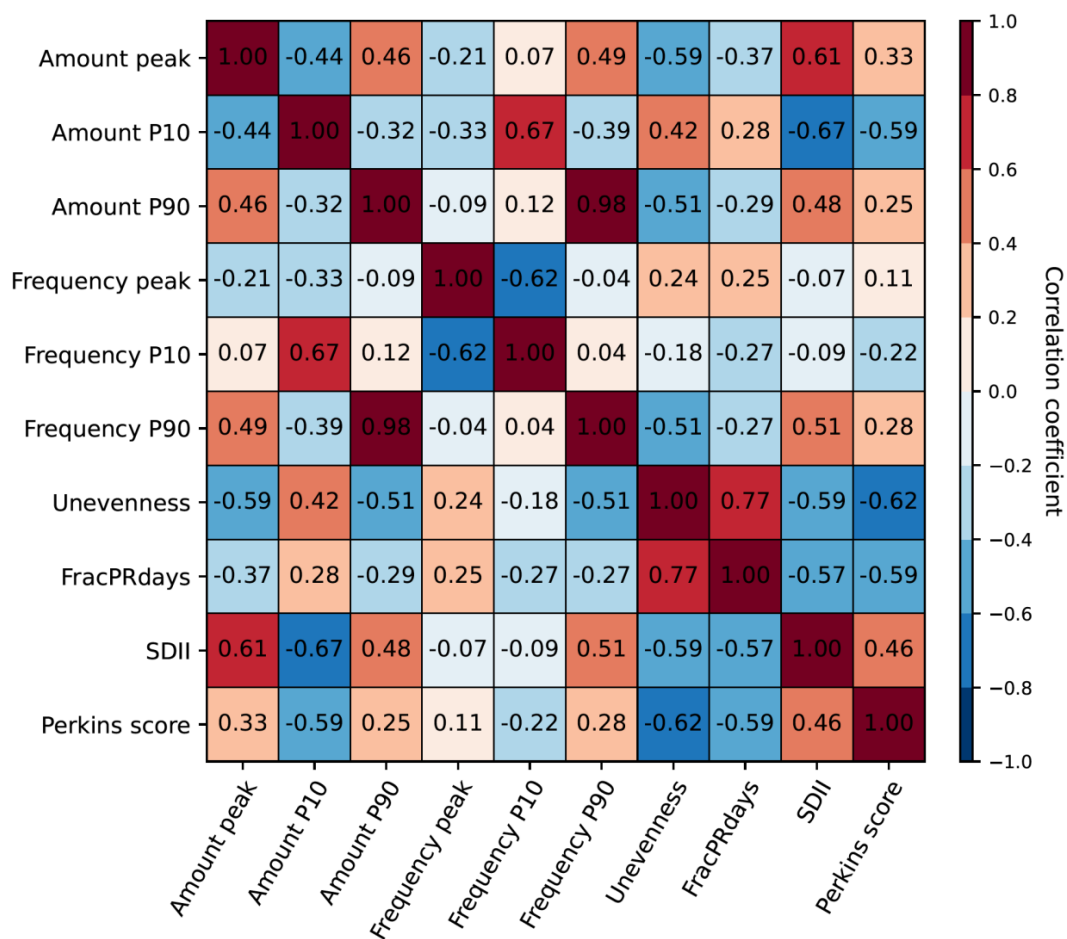


1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207

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.

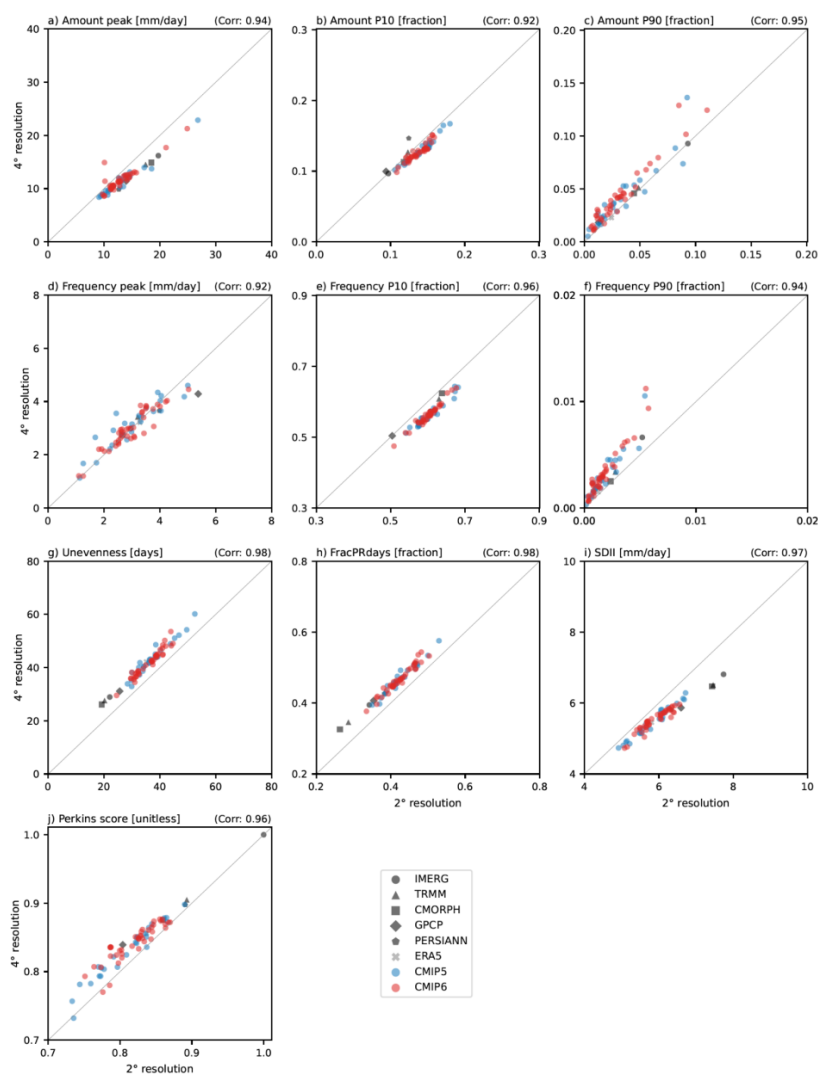


1208
 1209



1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222

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.



1223

1224

1225 Figure 14. Scatterplot between 2° and 4° interpolated horizontal resolutions in
 1226 evaluating precipitation distribution metrics for a) Amount peak, b) Amount P10, c)
 1227 Amount P90, d) Frequency peak, e) Frequency P10, f) Frequency P90, g) Unevenness,
 1228 h) FracPRdays, i) SDII, and j) Perkins score. The metric values are calculated for the
 1229 modified IPCC AR6 regions and then weighted averaged globally. Black, gray, blue, and
 1230 red marks indicate the satellite-based observations, reanalysis, CMIP5 models, and
 1231 CMIP6 modes, respectively. The number in the upper right of each panel is the
 1232 correlation coefficient between the metric values in 2° and 4° resolutions across all
 1233 observations and models.

1234