



1 **Evaluation of atmospheric rivers in reanalyses and climate models in a new**
2 **metrics framework**

3

4 Bo Dong¹, Paul Ullrich¹, Jiwoo Lee¹, Peter Gleckler¹, Kristin Chang¹, Travis A. O'Brien^{2,3}

5

6 1. Lawrence Livermore National Laboratory, Livermore, CA, USA

7 2. Department of Earth and Atmospheric Sciences, Indiana University,
8 Bloomington, IN, USA

9 3. Climate and Ecosystem Sciences Division, Lawrence Berkeley National Lab,
10 Berkeley, CA, USA

11

12

13 Correspondence to: Bo Dong (dong12@llnl.gov)

14

15

16 **Key points:**

17 1. A metrics package designed for easy analysis of AR characteristics and statistics is
18 presented

19 2. The tool is efficient for diagnosing systematic AR bias in climate models, and useful
20 for evaluating new AR characteristics in model simulations

21 3. In climate models, landfalling AR precipitation shows dry biases globally, and AR
22 tracks are farther poleward (equatorward) in the north and south Atlantic (south Pacific
23 and Indian Ocean)

24

25

26 **Abstract**

27

28 We present a suite of new atmospheric river (AR) metrics that are designed for quick
29 analysis of AR characteristics and statistics in gridded climate datasets such as model
30 output and reanalysis. This package is expected to be particularly useful for climate
31 model evaluation. The metrics include mean bias and spatial pattern correlation, which
32 are efficient for diagnosing systematic AR biases in climate models. For example, the
33 package identifies that in CMIP5 and CMIP6 models, AR tracks in the south Atlantic are
34 positioned farther poleward compared to the ERA5 reanalysis, while in the south
35 Pacific, tracks are generally biased towards the equator. For the landfalling AR peak
36 season, we find that most climate models simulate a completely opposite seasonal
37 cycle over western Africa. This tool is also useful for identifying and characterizing
38 structural differences among different AR detectors (ARDTs). For example, ARs
39 detected with the Mundhenk algorithm exhibit systematically larger size, width and
40 length compared to the TempestExtremes (TE) method. The AR metrics developed



41 from this work can be routinely applied for model benchmarking and during the
42 development cycle to trace performance evolution across model versions or generations
43 and set objective targets for the improvement of models. They can also be used by
44 operational centers to perform near real-time climate and extreme events impact
45 assessment as part of their forecast cycle.

46

47

48 **1. Introduction**

49

50 Atmospheric rivers (ARs) are dynamically driven, synoptic-scale filamentary structures
51 of water vapor jets that play important roles in the global water cycle and regional
52 weather and hydrology (Ralph et al. 2013; Gimeno et al. 2014; Shields et al. 2019;
53 Payne et al. 2020; O'Brien et al., 2022). These narrow, concentrated corridors of
54 moisture in the atmosphere can carry an immense amount of water, often compared to
55 the flow of multiple major rivers combined (Ralph and Dettinger, 2011), and account for
56 a substantial portion, more than 90% of the poleward water vapor transport (Zhu and
57 Newell, 1998; Newman et al. 2012; Ullrich et al. 2021). When approaching landmasses
58 or interacting with mountainous regions, ARs usually bring extreme weather inland,
59 such as heavy rainfall and strong wind, leading to severe flooding and landslides,
60 causing devastating damages to natural landscapes, agricultural fields, infrastructure,
61 human settlements, and disruption to businesses and services with significant economic
62 losses (Ralph et al., 2006; Leung and Qian, 2009; Neiman et al., 2011; Neiman et al.,
63 2013; Gershunov et al., 2017).

64

65 Previous studies have developed numerical algorithms for objective identification of ARs
66 (e.g., Neiman et al., 2009; Dettinger, 2011; Ralph et al., 2013; Mundhenk et al. 2016;
67 Ullrich and Zarzycki 2017; Ullrich et al., 2021). As noted by O'Brien et al. (2022), the
68 different choices made by ARDT developers essentially amount to different definitions
69 of ARs, all of which are qualitatively consistent with the definition in the AMS glossary
70 (Ralph et al., 2018). ARDTs are generally threshold-based, mostly using the intensity of
71 moisture transport with some geographical constraints that limit the AR spatial extent
72 and some geometrical constraints that preserve their nature as “long and narrow”
73 filaments of moisture. For example, the Mundhenk algorithm (Mundhenk et al. 2016)
74 calculates integrated water vapor transport (IVT) anomalies relative to the historical
75 period and uses a fixed relative threshold to identify ARs that are above a certain
76 percentile of the historical simulation. The TempestExtremes (TE; Ullrich et al. 2021)
77 method, as another example, uses relative threshold on the Laplacian of the IVT field
78 rather than the IVT field itself. Although AR detectors (ARDTs) are usually designed
79 with particular research questions in mind, they have widely facilitated broader studies



80 of AR characteristics and impacts (Shields et al., 2018; Rutz et al., 2019; O'Brien et al.,
81 2022).

82

83 The number of climate models under active development and used in the research
84 community has increased substantially in recent decades, with many supporting
85 multiple configurations and parameterization choices. Meanwhile, newer versions of
86 ARDTs have been developed , along with newer observational data products. As such,
87 routine evaluation of ARs during model development lifecycles requires a quantitative
88 climate data assessment evaluation workflow that is independent of ARDT and that
89 allows comparing AR characteristics from different ARDTs. We believe progress in
90 improving our understanding of ARs and their impacts could be accelerated with a
91 dedicated tool for calculating AR statistics in climate models and gridded data products.

92

93 Metrics have been widely used to quantify climate model performance in recent
94 decades (Taylor 2001; Gleckler et al. 2008; Wilks 2011; Zarzycki et al. 2021). Similarly,
95 a set of common metrics are also increasingly employed in AR studies over the past few
96 years, such as mean bias (Guan and Waliser 2017; Chapman et al. 2019), weighted
97 ensemble mean bias (Massoud et al. 2019), RMS error and relative RMS error (Guan
98 and Waliser 2017), spatial pattern correlation (Chapman et al. 2019; Huang et al. 2021),
99 ratio of spatial standard deviation (O'Brien et al. 2022), and skill scores for assessing
100 AR predictions (Wick et al. 2013, Nardi et al. 2018) and model performance (Zhang et
101 al. 2024). While these quantitative measures are case-specific and depend on the aim
102 of these studies, there is value in synthesizing commonly used metrics in one
103 comprehensive analysis tool.

104

105 In this paper, we propose a set of metrics that is designed for easy quantification of AR
106 characteristics and statistics in all types of gridded climate data, with the expectation
107 that such a metric suite would be useful for climate model evaluation. Following the
108 introduction, section 2 describes the general design of the AR metrics. Section 3
109 presents several example model evaluation applications of using the metrics evaluation
110 package. Conclusions and discussion are in section 4.

111

112

113 **2. Data and method**

114

115 **2.1 Input data**

116

117 The input data to the metrics package includes AR “tags” and optional climate variables
118 of interest that are concurrent with AR activities, such as precipitation, winds, and
119 temperature (Fig. 1). The AR tags can be products of any regional or global AR detector



120 (ARDT), including those based on relative (e.g., TempestExtremes or TE; Ullrich and
121 Zarzycki 2017; Ullrich et al. 2021), fixed-relative (e.g., Mundhenk_v3; Mundhenk et al.
122 2016), and absolute (e.g., Lora_v2; Lora et al. 2017) thresholds to the moisture field.

123

124 For applications in section 3, we run and compare the TE ARDT on the 6-hourly
125 integrated water vapor transport (IVT) data from three reanalysis products - ERA5
126 (Hersbach et al. 2020), MERRA-2 (Gelaro et al. 2017) and JRA-55C (Japan
127 Meteorological Agency, Japan 2015) to obtain AR tags for reanalyses. Given its longer
128 data record and finer model resolution, we subsequently use ERA5 as the default
129 reference in this study. To demonstrate how results are sensitive to the choice of
130 ARDTs, we then use the Mundhenk_v3 tags from ERA5 data.

131

132 To evaluate ARs in climate models, we use the archived AR tags from the Atmospheric
133 River Tracking Method Intercomparison Project (ARTMIP) Tier 2 experiment, which is
134 based on the coupled CMIP model simulations for the historical and 21st century
135 projection periods. (Shield et al. 2019, Rutz et al 2019, O'Brien et al, 2022). The tag
136 data include six of the CMIP5 models (CCSM4, CSIRO-Mk3-6, CanESM2, IPSL-CM5A-
137 LR, IPSL- CM5B-L, and NorESM1-M) and 3 of the CMIP6 models (BCC-CSM2-MR,
138 IPSL-CM6A-LR, MRI-ESM2-0). For model evaluation purposes in our application
139 examples, only TE tags from the archive are selected.

140

141 We further use simulations from the Energy Exascale Earth System Model (E3SM;
142 Golaz et al. 2019, Caldwell et al. 2019) high resolution (HR, 0.25°, ~28 km grid) and low
143 resolution (LR, 1°, ~111 km grid) experiments to examine the sensitivity of ARs to
144 model resolution. Except for their different horizontal grid spacing, both E3SM-HR and
145 E3SM-LR use an identical set of physical parameters, and the simulations follow a
146 similar protocol of the Coupled Model Intercomparison Project Phase 6 (CMIP6; Eyring
147 et al. 2016).

148

149 For the evaluation of AR characteristics, statistics gauging the consistency of latitude,
150 longitude, width, length, and size are required as the input for metrics. In our case, we
151 use the 'BlobStats' tool (Ullrich et al. 2021) to calculate the statistics, where latitude and
152 longitude are weighted by the moisture field, width and length are based on principle
153 component analysis (PCA; Inda-Díaz et al. 2021), and size is based on a count of the
154 number of contiguous grid cells in the feature. This tool can be called and run within the
155 AR metrics workflow, with a separate installation. Users can also optionally use their
156 preferred statistical package for AR geometry calculation and then feed the data back to
157 the metrics workflow.

158



159 2.2 Geographical Regions

160

161 In this tool package, the AR metrics are calculated based on the data in user-defined
162 geographic domains. In Fig. 1, the upper right panel shows examples of regions that
163 were selected for landfalling AR diagnostics (red boxes in the panel, lat-lon boundaries
164 are listed in the supplementary table S1). These regions, mostly located in the west
165 coast of continents, are known to have frequently observed AR landfalls (Guan and
166 Waliser 2015, Algarra et al. 2020). We purposely use rectangular region boundaries for
167 easy use of the metrics tool, such that rather than needing a regional mask file, users
168 can quickly sub-select the data by declaring latitude and longitude bounds of any
169 specific region. For AR statistics, we group global ARs into 5 major ocean basins – the
170 North Pacific, South Pacific, North Atlantic, South Atlantic, and South Indian Ocean
171 (blue boxes in Fig. 1 upper right panel; lat-lon coordinates in table S1 in the
172 supplement).

173

174 2.3 Metrics

175

176 2.3.1 Mean bias

177

178 We use mean bias to measure how close a climate data product is with respect to the
179 reference data, calculated as

$$180 \quad \bar{b} = \bar{x} - \bar{y}$$

181 where \bar{x} is the arithmetic mean of the test variable x with sample size n , given by

182

$$183 \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

184

185 and similarly, the \bar{y} is the arithmetic mean of the reference variable

186

187 The statistical significance of the mean bias is measured using the Z-test, with the test
188 statistics (z-score) formulated as

189

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - (\bar{\mu}_1 - \bar{\mu}_2)}{\sqrt{\frac{s_1^2}{n_1-1} + \frac{s_2^2}{n_2-1}}}$$

190

191

192 where \bar{x}_i is sample arithmetic mean, μ_i is population mean, s_i is sample variance, and
193 n_i is sample size. A positive z-score indicates that the value is above the mean. The
194 higher the z-score, the further above the mean the value is, and vice versa. A result is



195 considered statistically significant at the 95% confidence level if the magnitude of the z-
196 score is greater than 1.96.

197

198 When comparing across different variables, a commonly used measure is the
199 normalized bias, with the data normalized by the standard deviation of the reference
200 field. In this study, we simply use z-score as the normalized bias, as it incorporates both
201 bias and statistical significance in one succinct formula.

202

203 2.3.2 Spatial pattern similarity

204

205 The spatial pattern correlation is a measure used to quantify the similarity between two
206 spatial fields without reflecting the magnitude of the difference. Here we compute the
207 spatial pattern correlation using the Pearson correlation coefficient:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

208

209 where, x_i and y_i are the values of the two spatial patterns at location i (or grid point i in
210 gridded data product), \bar{x} and \bar{y} are the means of the values of the two patterns, and n is
211 the total number of locations. This equation essentially measures the degree to which
212 the values of the two spatial patterns vary together. If they vary together perfectly, r will
213 be 1. If they vary together inversely, r will be -1. If there's no linear relationship between
214 the patterns, r will be 0.

215

216 The statistical significance of correlation is determined by the two-tailed p-value of the
217 cumulative distribution function (CDF) of the t-statistic, as

218

$$219 \quad p = 2 \times (1 - \text{CDF}(t))$$

220

221 The the t-statistic t is given by

$$t = r \times \frac{\sqrt{n_e}}{\sqrt{1 - r^2}}$$

222

223 where r is the correlation coefficient, and n_e is the effective sample size. Although there
224 are a number of methods to estimate the effective geographic sample size (e.g., Griffith
225 2013), given that ARs present notable seasonal and interannual latitudinal shift
226 patterns, we propose a new method to estimate n_e as the number of Principal
227 Component Analysis (PCA) modes required to explain more than 95% of the total



228 variance in the AR tag data. The cumulative variance explained by the principal
229 components is expressed as

$$n_e = \min \left\{ n_e \mid \frac{\sum_{i=1}^{n_e} \lambda_i}{\sum_{i=1}^p \lambda_i} > 0.95 \right\}$$

230
231

232 where the λ_i are the eigenvalues of the spatial correlation matrix of the data, and p is the
233 total number of principal components. Estimating n_e based on ERA5 reanalysis data, we
234 find that the effective sample sizes for spatial pattern correlation are generally small,
235 ranging from 14 - 27 for the 5 ocean basins (Table S2 in supplementary information).

236

237 2.3.3 Temporal detection similarity

238

239 The AR binary occurrence time series refers to a binary time series equal to one when
240 an AR is present in a given region and zero otherwise. The overlap between two AR
241 occurrence time series is measured by the Intersection over Union (IoU) metric. The
242 metric is written as

$$IoU(A, B) = \frac{\sum |A \cap B|}{\sum |A \cup B|}$$

243

244 where, A and B are binary AR occurrence time series. The IoU is useful for gauging the
245 degree of temporal similarity of ARs detected in different ARDTs.

246

247 2.3.4 Metrics and diagnostics implementation

248

249 The metrics and diagnostics are pre-defined in the metrics framework, and they are fully
250 customizable. Table 1 lists all the AR metrics and diagnostics used in this study. The
251 AR metrics are composed of AR properties (as shown in the top row) and evaluation
252 metrics. Similarly, the AR diagnostics are composed of AR properties and statistical
253 diagnostics. The number of regions that these metrics are applied to are indicated by
254 the numbers in the table. The metrics code is python-based, and it handles gridded AR
255 tag and climate data using xCDAT (Xarray Climate Data Analysis Tools,
256 <https://xcdat.readthedocs.io>), which is an extension of the Xarray package
257 (<https://xarray.pydata.org>).

258

259

260 3. Metrics applications

261

262 In this section, we present five example applications using the metrics tool for assessing
263 ARs in climate models, including evaluation of AR frequency and characteristics,



264 comparison of ARs in high- and low-resolution simulations, sensitivity of ARs to choice
265 of ARDT, precipitation bias associated with ARs and landfalling AR seasonality.

266

267 **3.1 AR characteristics in CMIP5 and CMIP6 models**

268

269 **3.1.1 AR frequency**

270

271 We first analyze the pattern of AR occurrence frequency over a 10-year period (1979-
272 1988) for the five major ocean basins from section 2.2. From the spatial distribution of
273 the AR frequency, we calculate the pattern correlation between selected climate models
274 and ERA5. The spatial pattern correlation coefficient is shown in Fig. 2. Notably the
275 correlations are statistically significant for all models and regions. This suggests that
276 climatologically, climate models simulate AR density and spatial distribution that broadly
277 resemble reanalysis on planetary scale. This is evidenced in the spatial AR occurrence
278 density maps in Fig. 3 (a-b) and (d-e).

279

280 The high spatial correlation is mainly a result of the similar spatial gradient of the AR
281 frequencies, rather than the similar magnitude of the frequency at each grid point in two
282 datasets. For instance, if the AR frequency values in one map are doubled compared to
283 those on the other map, the spatial patterns, or spatial structures of the two, can still be
284 perfectly correlated. Since climatologically ARs are largely clustered along the storm
285 track, with nearly no presence over a large portion of the basin domain, it is natural that
286 the pattern correlations are significant in most cases. Similar high pattern correlations of
287 AR frequencies are also noted in other studies (e.g., Huang et al. 2020; Guan et al.
288 2023). In other words, the spatial correlation coefficient is not that indicative for the
289 magnitude resemblance of the AR spatial frequency. Therefore, these metric results can
290 be better interpreted together with AR frequency maps.

291

292 While the spatial correlation coefficient synthesizes the level of pattern consistency,
293 difference maps further reveal the spatial discrepancies. For example, Fig. 3c shows
294 that South Pacific AR tracks shift farther towards the equator in the CSIRO model than
295 in ERA5. While in the North Atlantic basin (Fig. 3f), AR tracks are displaced more
296 poleward in the BCC model. The further north AR location is likely associated with the
297 poleward jet stream bias in CMIP6 models (Bracegirdle et al. 2020; Harvey et al. 2020).

298

299 **3.1.2 AR geometric features in major ocean basins**

300

301 The portrait plots in Fig. 4 show normalized biases (as z-score) of AR characteristics in
302 climate models for the 5 major ocean basins. Several striking results emerge. For
303 instance, in the North Pacific, the CMIP5 and CMIP6 AR geometry, in terms of width
304 and length, are significantly smaller than the ERA5 reanalysis. One possible cause of



305 such bias is that the AR blobs detected with TE in the relatively lower resolution climate
306 models are geometrically less curvy, and less pointy at the ends. Fig. S1 shows an
307 example time slice of AR blobs in the ERA5 and BCC model. It is clear that the
308 highlighted AR blob in the BCC model exhibits a “cut-off” feature at both ends, thus
309 shorter in length than the ERA5 reanalysis. And although visually the blob is wider, the
310 PCA based width is actually narrower due to its less curvy blob geometry. In contrast,
311 for all other ocean basins, the AR sizes (area) are generally bigger in climate models.
312 The figures also show notable latitudinal model AR biases, such that compared to the
313 reanalysis, ARs tend to shift towards higher latitudes in the North and South Atlantic
314 and biased towards the equator in the South Pacific and Indian Ocean.

315

316 Fig. 4 also helps identify outliers of a specific model or variable. For example, although
317 most climate models tend to simulate larger ARs than observed (indicated by the
318 positive values in the area columns), one notable exception is the CanESM2 model
319 which has significantly smaller AR width, length, and area than other models and ERA5
320 reanalysis. Taking a closer look into the AR width and length in the North Pacific in Fig.
321 5, we see that CanESM2 simulates more smaller ARs and fewer bigger ARs than the
322 reanalysis, resulting in negative mean biases. This type of histogram helps us better
323 understand the AR distribution discrepancies.

324

325 Another example is from the CCSM4 model simulations. The higher bounds of the
326 model histogram in nearly all fields indicate that the CCSM4 model simulates more ARs
327 than the reanalysis, with bigger size indicated as taller area bars in Fig. 5c. The higher
328 ARs counts in the model are mostly located in the high latitudes and the tropics south of
329 20°N (Fig. 5a), spreading across all longitude (Fig. 5b). Fig. 5d and 5e show that the
330 additional ARs in CCSM4 are narrower and/or longer in shape.

331

332 **3.2 ARs in high and low resolution E3SM simulations**

333

334 We now apply the metrics and diagnostics identified in section 2.3.4, including the mean
335 bias of AR latitude, longitude, area, width and length over 5 ocean basins, and AR
336 induced precipitation over 16 landfall regions, to evaluate and compare AR
337 characteristics in the E3SM HR and LR simulations. ARs in both HR and LR exhibit
338 similar structural differences compared to the ERA5 (Fig. 6a, b). They are bigger in
339 terms of area, width, and length, and biased towards higher latitudes in the North Pacific
340 and South Atlantic. Zonally, ARs in E3SM are more westward distributed in the North
341 Pacific, and more eastward distributed in the North Atlantic and South Pacific. One
342 difference we see between the two experiments is that in the North Atlantic basin, AR
343 tracks in the HR are shifted more northward than in the LR simulation.

344



345 Figure 6c shows AR differences between E3SM HR and LR models. The most
346 noticeable differences are that the HR simulates wider and longer ARs than the LR
347 model over all ocean basins. The AR size, in the area column, however, shows mixed
348 results which are not consistent with systematic biases in width and length. This is
349 probably because of different AR geometric properties in the HR and LR simulations.
350 For example, in Supplementary Figure S2, the highlighted AR blob in the North Atlantic
351 is longer but smaller in the LR compared to the one in the HR simulation. Latitudinally,
352 AR distributions show hemispheric contrast, as compared to the LR, ARs in HR are
353 located more southward in the Pacific sector but more northward in the Atlantic sector.

354

355 Figure 7 shows AR characteristic distribution in the North Pacific for E3SM HR, LR and
356 ERA5. Apparently, E3SM produces more AR events than the reanalysis in nearly all
357 fields and across all scales. We also evaluated the precipitation associated with
358 landfalling ARs in California in both HR and LR simulations, as in Fig. 8. It is notable
359 that both models simulate systematically higher precipitation than ERA5 for all rainfall
360 intensity categories. It is also clear that the precipitation bias in HR simulation is larger
361 than LR simulation, except in the light rainfall ($< \sim 6\text{mm/day}$) category. Similarly, better
362 topographic representation in high resolution version of the model does not improve
363 precipitation simulation is also reported in Harrop et al. (2023), especially when the bias
364 in the low resolution model is substantially high.

365

366 **3.3. Sensitivity of AR characteristics to ARDT**

367

368 In this application of the metrics package, we examine how ARs in ERA5 are sensitive
369 to the choice of ARDT. In addition to TE-based AR tags, we use AR tags detected using
370 the Mundhenk_v3 algorithm for comparison. Despite significant differences in their
371 associated algorithms, results from ARTMIP showed their performance was similar and
372 close to the mean among all ARDTs (Shields et al., 2018). Table 2 shows agreement of
373 landfalling ARs detected using these two ARDTs, as % values of IoU (AR concurrence
374 normalized by total occurrence of the ARs in both methods). The level of consistency
375 ranges from 56% to 83%, which suggests that TE and Mundhenk detect ARs
376 concurrently most of the time, but with asynchronous discrepancies, possibly at the
377 timing of the landfall and the end of the AR life cycle.

378

379 For AR characteristics over the oceans, the Mundhenk method detects larger ARs in
380 area, width, and length compared to TE (Fig. 9). ARs are also present at more
381 northward latitudes with Mundhenk than TE. Zonally, AR distributions exhibit more
382 hemispherical contrast, with Mundhenk showing more westward located ARs in the
383 Pacific sector but more eastward located ARs in the Atlantic sector.

384



385 **3.4 Landfalling AR precipitation in CMIP5/6 models**

386

387 Precipitation is an important indicator of the intensity of a landfalling AR. Here we
388 evaluate landfalling AR precipitation in the CMIP5 and CMIP6 models, with the ERA5
389 reanalysis and MSWEP (Beck et al. 2017) gridded product as reference. Fig. 10 shows
390 that compared to the observations, landfalling precipitation differences in the models are
391 generally much larger than in reanalysis. The models show dry biases in most regions,
392 particularly large in California, Pacific Northwest, Iceland and Greenland.

393

394 As it is unclear if these biases are mainly due to general precipitation biases, or AR
395 activity bias, we further examine model precipitation bias diagnostics regardless of AR
396 activity (Fig. 11a) and AR frequency bias metrics (Fig. 11b) separately. For total
397 precipitation in the models, structural biases as in Fig. 10 are absent, but AR landfalls
398 are less frequent in the Pacific Northwest, Iceland, and Greenland. This suggests that
399 the systematic dry AR precipitation biases over these regions are primarily due to the
400 insufficient number of landfalling ARs in the models. For California, similar results do not
401 hold for all the models, for example, total precipitation in CCSM4 is higher than the
402 reanalysis and AR landfalls are more frequent, but the AR-related rainfall has a
403 significant dry bias. This suggests that landfalling ARs in CCSM4 are less intense,
404 suggesting a potential direction for model improvement.

405

406 **3.5 Landfalling AR peak day**

407

408 **3.5.1 Comparison among reanalyses**

409

410 Seasonality of AR landfalls is one of the important metrics for understanding AR
411 variability and impacts. Here we analyze landfalling AR seasonality over various regions
412 of the globe among three reanalysis products. We perform a Fourier transform on the
413 10-year long-term daily mean AR histogram to find its peak date based on the phase of
414 the first Fourier mode. Results indicate that the AR peak days agree well among
415 reanalyses for most regions, with small differences of only a few days. Large
416 discrepancies are noted for Australia and western Africa: In Australia, AR landfall peaks
417 nearly a month behind in JRA-55C than MERRA-2, while in west Africa, AR landfall in
418 MERRA-2 peaks 46 days behind ERA5.

419

420 Details of these differences are depicted in the histogram plots. For West Africa, AR
421 landfalls have two peaks in ERA5 and MERRA-2, one being in September, followed by
422 another peak in November. In ERA5, the peak in November is the main peak, while in
423 MERRA-2, the September peak is comparable to the November peak, resulting in an
424 earlier peak day from the Fourier phase spectrum. JRA-55C, in contrast, has only one



425 peak in November, and the AR landfall event counts are fewer than the other two
426 products over the entire year, indicative of smaller year to year variability.

427

428 Seasonal distribution of AR landfalls in Australia in the three reanalyses exhibit similar
429 differences to those in western Africa. In ERA5 and MERRA-2, there are two peaks in
430 February and June, but only one peak presents in JRA-55C in June. This explains the
431 relative late peak day in JRA-55C. While the main peak in ERA5 is in June, in MERRA-
432 2, the main peak is in February, which is consistent with the metrics result that MERRA-
433 2 has the earliest peak day. Similarly, the JRA-55C has a smaller number of landfalling
434 ARs, although the interannual variability is comparable to the other two reanalyses.

435

436 3.5.2 Evaluation of climate models

437

438 Figure 13 shows CMIP5 and CMIP6 models' performance in simulating AR peak
439 season compared to ERA5 reanalysis. To explore how model biases compare to the
440 discrepancies among reanalyses, we also include AR peak day bias for MERRA-2 and
441 JRA-55C reanalysis in the left two columns of the metrics plot. Perhaps unsurprisingly,
442 the model spread is much larger than the spread among reanalysis products, which are
443 tightly constrained by data assimilation.

444

445 In regions like South America, Baja, UK and Western Europe, the models show
446 systematic late peak biases, and in South Africa, AR peaks earlier than the reanalyses.
447 The exact cause of these structural biases in the models is likely indicative of persistent
448 and ubiquitous timing issues in the shift of the storm track that is common among
449 models. It is worth noting that the model biases in the West Africa region are
450 significantly larger than other regions, with peak day difference up to 6 months as
451 compared to the reanalysis. Looking at the AR counts histograms over the course of the
452 year in this region in the CCSM3 and MRI-ESM2-0 models (Fig. 14), it is clear that AR
453 landfall seasonality in both models is completely out of phase with ERA5. This is
454 especially true for the MRI-ESM2-0 model, where AR landfall peaks in June, which is in
455 opposition to the climatology in ERA5. The large discrepancy is probably because of the
456 large spread in the atmospheric circulations in climate models over the West Africa
457 region, as large spread among CMIP5/6 models in capturing atmospheric dynamic
458 responses (Monerie et al. 2020), the lack of jet-rainfall coupling (Whittleston et al. 2017),
459 and bias in simulating mesoscale convective systems (Jenkins et al. 2002) in climate
460 models are noted. Although high resolution regional modeling may be capable of
461 improving rainfall in this region (Sylla et al. 2009), the dynamics-rainfall coupling does
462 not appear to be improved in high resolution global models such as the E3SM (Caldwell
463 et al. 2019; Golaz et al. 2019). Therefore, challenges remain in modeling the AR water
464 cycle in west Africa.



465

466 **4. Summary and discussion**

467

468 In this study we have introduced a workflow for the objective evaluation of ARs in
469 climate models and reanalysis, and have illustrated the potential for its use with five
470 example case-studies to illustrate the scope of potential applications. The metrics-based
471 analyses are designed for systematic diagnosis of AR biases in climate models. For
472 example, applying the package to CMIP5 and CMIP6 models, we have shown that AR
473 tracks in the south Atlantic are positioned farther poleward compared to the ERA5
474 reanalysis, while in the south Pacific, tracks are biased towards the equator. Over
475 western Africa, we found that most climate models do a poor job at capturing the AR
476 peak season. In addition to model evaluation, we have shown how our tool can be used
477 to identify structural differences resulting from the choice of AR detector (ARDT). For
478 instance, we demonstrated that ARs detected with the Mundhenk method are
479 systematically larger in size, width and length compared to TE.

480

481 The workflow and metrics presented in this study can be used for a variety applications,
482 e.g., to contrast the differences between AR features in historical and future scenarios
483 as simulated by climate models. Objectively quantifying projected changes in landfall
484 frequency, duration, and intervals between landfall events are of particular interest.
485 Further confidence in this and other model evaluation applications can be gained by
486 assessing what impact the choice of the ARDT can have on any conclusions concerning
487 model quality. Our tool makes this and other sensitivity tests more tractable.

488

489 Our tool also pools a diverse suite of established and newly introduced AR metrics into
490 one framework, facilitating objective evaluation of ARs with a diverse suite of input data,
491 as well as intercomparison of ARs as simulated by multiple climate models. These
492 metrics can be routinely applied for model benchmarking and during development
493 cycles to monitor changes in AR characteristics across model versions or generations
494 and set objective targets for the improvement of models. One expected application is
495 the routine benchmarking of AR in simulations with increasingly higher resolution
496 models. More frequent metrics evaluation of simulated ARs such as this could further
497 our understanding of model bias and error characteristics, and potentially assist
498 developers in making choices associated with new model versions. Furthermore, it
499 effectively provides a quantitative measure for operational centres to perform near real-
500 time climate and extreme events impact assessment along with their forecast cycles,
501 which can facilitate their decision-making process.

502

503 Our metrics tool is developed with Xarray (Hoyer et al., 2017), XCDAT (Vo et al., 2024),
504 and the PCMDI Metrics Package (PMP; Lee et al. 2024), which are compatible with one



505 another, readily available and easy to install. At the time of the submission of this
506 manuscript, our tool is being configured to be a part of the PMP. Looking forward, we
507 welcome community contributions to successive development of the package. Inspired
508 by Zarzycki et al. (2021), there is also a potential that these metrics can be applied for
509 research beyond ARs, such as mesoscale meteorological features, regional
510 hydrological extremes such as floods and droughts, and large-scale climate modes.

511
512
513
514

515 **Acknowledgment**

516

517 This work is performed under the auspices of the U.S. Department of Energy (DOE) by
518 Lawrence Livermore National Laboratory (LLNL) under Contract No. DE-AC52-
519 07NA27344 and is mainly supported by the Regional and Global Model Analysis
520 (RGMA) program of the U.S. DOE Office of Science (OS) Biological and Environmental
521 Research (BER) program. This material is based upon work supported by the U.S.
522 Department of Energy, Office of Science, Office of Biological and Environmental
523 Research, Climate and Environmental Sciences Division, Regional & Global Model
524 Analysis Program. Resources of the National Energy Research Scientific Computing
525 Center (NERSC) were used. The research was partially supported by the Office of
526 Science of the U.S. Department of Energy under Contract Number DE-AC02-
527 05CH11231 and under award Number DE-SC0023519. This research was also
528 supported in part by the Environmental Resilience Institute, funded by Indiana
529 University's Prepared for Environmental Change Grand Challenge initiative and in part
530 by Lilly Endowment, Inc., through its support for the Indiana University Pervasive
531 Technology Institute. We acknowledge the World Climate Research Programme, which,
532 through its Working Group on Coupled Modeling, coordinated and promoted CMIP6.
533 We thank the climate modeling groups for producing and making available their model
534 output, the Earth System Grid Federation (ESGF) for archiving the data and providing
535 access, and the multiple funding agencies that support CMIP6 and ESGF. The authors
536 acknowledge Antony Hoang and Ana Ordonez for computing and technical support, and
537 Christine Shields, Yang Zhou and Allison Collow for their help on data and discussion.

538

539 **Code and data availability**

540 The metrics framework code is available on github. Users are recommended to install
541 full PMP package (http://pcmdi.github.io/pcmdi_metrics/install.html) in order to have the
542 environment and python packages to run the metrics code.

543

544 **Author contribution**



545 BD implemented the codes and developed the diagnostic results. All authors
546 contributed to the writing of the manuscript.

547

548 **Competing interests**

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

551

552 **References**

553

554 Algarra, I., Nieto, R., Ramos, A. M., Eiras-Barca, J., Trigo, R. M., & Gimeno, L. (2020).
555 Significant increase of global anomalous moisture uptake feeding land-falling
556 atmospheric rivers. *Nature communications*, 11 (1), 5082.

557 Beck, H. E., Van Dijk, A. I., Levizzani, V., Schellekens, J., Miralles, D. G., Martens, B., &
558 De Roo, A. (2017). MSWEP: 3-hourly 0.25 global gridded precipitation (1979–2015)
559 by merging gauge, satellite, and reanalysis data. *Hydrology and Earth System
560 Sciences*, 21(1), 589-615.

561 Caldwell, P. M., Mamejtanov, A., Tang, Q., Van Roekel, L. P., Golaz, J. C., Lin, W., ... &
562 Zhou, T. (2019). The DOE E3SM coupled model version 1: Description and results
563 at high resolution. *Journal of Advances in Modeling Earth Systems*, 11(12), 4095-
564 4146.

565 Chapman, W. E., Subramanian, A. C., Delle Monache, L., Xie, S. P., & Ralph, F. M.
566 (2019). Improving atmospheric river forecasts with machine learning. *Geophysical
567 Research Letters*, 46(17-18), 10627-10635.

568 DeFlorio, M. J., Waliser, D. E., Guan, B., Lavers, D. A., Ralph, F. M., & Vitart, F. (2018).
569 Global assessment of atmospheric river prediction skill. *Journal of
570 Hydrometeorology*, 19(2), 409-426.

571 Dettinger, M. D., Ralph, F. M., Das, T., Neiman, P. J., & Cayan, D. R. (2011). Atmospheric
572 rivers, floods and the water resources of california. *Water*, 3 (2), 445–478.

573 Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., & Taylor, K.
574 E. (2016). Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6)
575 experimental design and organization. *Geoscientific Model Development*, 9(5),
576 1937-1958.

577 Gelaro, R., McCarty, W., Suárez, M. J., Todling, R., Molod, A., Takacs, L., . . . others
578 (2017). The modern-era retrospective analysis for research and applications,
579 version 2 (merra-2). *Journal of climate*, 30 (14), 5419–5454.

580 Gershunov, A., Shulgina, T., Clemesha, R. E., Guirguis, K., Pierce, D. W., Dettinger, M.
581 D., . . . others (2019). Precipitation regime change in western north america: the role
582 of atmospheric rivers. *Scientific reports*, 9 (1), 9944.

583 Gimeno, L., Nieto, R., Vázquez, M., & Lavers, D. A. (2014). Atmospheric rivers: Amini-
584 review. *Frontiers in Earth Science*, 2 , 2.



- 585 Gleckler, P. J., Taylor, K. E., & Doutriaux, C. (2008). Performance metrics for climate
586 models. *Journal of Geophysical Research: Atmospheres*, 113(D6).
- 587 Golaz, J. C., Caldwell, P. M., Van Roekel, L. P., Petersen, M. R., Tang, Q., Wolfe, J. D., ...
588 & Zhu, Q. (2019). The DOE E3SM coupled model version 1: Overview and
589 evaluation at standard resolution. *Journal of Advances in Modeling Earth Systems*,
590 11(7), 2089-2129.
- 591 Griffith, D. A. (2013). Establishing qualitative geographic sample size in the presence of
592 spatial autocorrelation. *Annals of the Association of American Geographers*, 103(5),
593 1107-1122.
- 594 Guan, B., Molotch, N. P., Waliser, D. E., Fetzer, E. J., & Neiman, P. J. (2010). Extreme
595 snowfall events linked to atmospheric rivers and surface air temperature via satellite
596 measurements. *Geophysical Research Letters*, 37 (20).
- 597 Guan, B., & Waliser, D. E. (2017). Atmospheric rivers in 20 year weather and climate
598 simulations: A multimodel, global evaluation. *Journal of Geophysical Research:*
599 *Atmospheres*, 122(11), 5556-5581.
- 600 Guan, B., Waliser, D. E., & Ralph, F. M. (2023). Global application of the atmospheric
601 river scale. *Journal of Geophysical Research: Atmospheres*, 128(3),
602 e2022JD037180.
- 603 Harrop, B., Leung, L. and Ullrich P. (2023). Improving Simulations of Atmospheric Rivers
604 and Heat Waves in the Coupled E3SM. FY2023 First Quarter Performance Metric.
605 DOE/SC-CM-23-001
- 606 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., . . .
607 others (2020). The ERA5 global reanalysis. *Quarterly Journal of the Royal*
608 *Meteorological Society*, 146 (730), 1999–2049.
- 609 Hoyer, S. & Hamman, J., (2017). xarray: N-D labeled Arrays and Datasets in Python.
610 *Journal of Open Research Software*. 5(1), p.10. DOI:
611 <https://doi.org/10.5334/jors.148>
- 612 Huang, X., Swain, D. L., & Hall, A. D. (2020). Future precipitation increase from very high
613 resolution ensemble downscaling of extreme atmospheric river storms in California.
614 *Science advances*, 6(29), eaba1323.
- 615 Huang, J., Zhang, C., & Prospero, J. M. (2009). African aerosol and large-scale
616 precipitation variability over West Africa. *Environmental Research Letters*, 4(1),
617 015006.
- 618 Hui, W. J., Cook, B. I., Ravi, S., Fuentes, J. D., & D'Odorico, P. (2008). Dust-rainfall
619 feedbacks in the West African Sahel. *Water Resources Research*, 44(5).
- 620 Inda-Díaz, H. A., O'Brien, T. A., Zhou, Y., & Collins, W. D. (2021). Constraining and
621 characterizing the size of atmospheric rivers: A perspective independent from the
622 detection algorithm. *Journal of Geophysical Research: Atmospheres*, 126(16),
623 e2020JD033746.



- 624 Kobayashi, S., Ota, Y., Harada, Y., Ebita, A., Moriya, M., Onoda, H., . . . others (2015).
625 The JRA-55 reanalysis: General specifications and basic characteristics. *Journal of*
626 *the Meteorological Society of Japan. Ser. II* , 93 (1), 5–48.
- 627 Lee, J., P. J. Gleckler, M.-S. Ahn, A. Ordonez, P. Ullrich, K. R. Sperber, K. E. Taylor, Y.
628 Y. Planton, E. Guilyardi, P. Durack, C. Bonfils, M. D. Zelinka, L.-W. Chao, B. Dong,
629 C. Doutriaux, C. Zhang, T. Vo, J. Boutte, M. F. Wehner, A. G. Pendergrass, D. Kim,
630 Z. Xue, A. T. Wittenberg, and J. Krasting, (2024): Systematic and Objective
631 Evaluation of Earth System Models: PCMDI Metrics Package (PMP) version 3.
632 *Geoscientific Model Development*, 17, 3919–3948, doi: 10.5194/gmd-17-3919-
633 2024.
- 634
- 635 Leung, L. R., & Qian, Y. (2009). Atmospheric rivers induced heavy precipitation and
636 flooding in the western us simulated by the WRF regional climate model.
637 *Geophysical research letters*, 36 (3).
- 638 Lora, J. M., Mitchell, J. L., Risi, C., & Tripathi, A. E. (2017). North pacific atmospheric rivers
639 and their influence on western north America at the last glacial maximum.
640 *Geophysical Research Letters*, 44 (2), 1051–1059.
- 641 Massoud, E. C., Espinoza, V., Guan, B., & Waliser, D. E. (2019). Global Climate Model
642 Ensemble Approaches for Future Projections of Atmospheric Rivers. *Earth's Future*,
643 7: 1136–11511151.
- 644 Mundhenk, B. D., Barnes, E. A., & Maloney, E. D. (2016). All-season climatology and
645 variability of atmospheric river frequencies over the north pacific. *Journal of Climate*,
646 29 (13), 4885–4903.
- 647 Nardi, K. M., Barnes, E. A., & Ralph, F. M. (2018). Assessment of numerical weather
648 prediction model reforecasts of the occurrence, intensity, and location of
649 atmospheric rivers along the West Coast of North America. *Monthly Weather*
650 *Review*, 146(10), 3343–3362.
- 651 Neiman, P. J., Ralph, F. M., Moore, B. J., Hughes, M., Mahoney, K. M., Cordeira, J. M.,
652 & Dettinger, M. D. (2013). The landfall and inland penetration of a flood-producing
653 atmospheric river in arizona. part i: Observed synoptic-scale, orographic, and
654 hydrometeorological characteristics. *Journal of Hydrometeorology*, 14 (2), 460–484.
- 655 Neiman, P. J., Schick, L. J., Ralph, F. M., Hughes, M., & Wick, G. A. (2011). Flooding in
656 western washington: The connection to atmospheric rivers. *Journal of*
657 *Hydrometeorology*, 12 (6), 1337–1358.
- 658 Neiman, P. J., White, A. B., Ralph, F. M., Gottas, D. J., & Gutman, S. I. (2009). A water
659 vapour flux tool for precipitation forecasting. In *Proceedings of the institution of civil*
660 *engineers-water management (Vol. 162, pp. 83–94)*.
- 661 Newman, M., Kiladis, G. N., Weickmann, K. M., Ralph, F. M., & Sardeshmukh, P. D.
662 (2012). Relative contributions of synoptic and low-frequency eddies to time-mean



- 663 atmospheric moisture transport, including the role of atmospheric rivers. *Journal of*
664 *climate*, 25(21), 7341-7361.
- 665 O'Brien, T., Wehner, M., Payne, A., Shields, C., Rutz, J., Leung, L.-R., . . . others (2022).
666 Increases in Future AR Count and Size: Overview of the ARTMIP Tier 2 CMIP5/6
667 Experiment. *Journal of Geophysical Research: Atmospheres*, 127 (6).
- 668 Payne, A. E., Demory, M. E., Leung, L. R., Ramos, A. M., Shields, C. A., Rutz, J. J., ... &
669 Ralph, F. M. (2020). Responses and impacts of atmospheric rivers to climate
670 change. *Nature Reviews Earth & Environment*, 1(3), 143-157.
- 671 Ramachandran, J., & Aschheim, M. A. (2005). Sample size and error in the determination
672 of mode shapes by principal components analysis. *Engineering structures*, 27(14),
673 1951-1967.
- 674 Ralph, F., Coleman, T., Neiman, P., Zamora, R., & Dettinger, M. (2013). Observed
675 impacts of duration and seasonality of atmospheric-river landfalls on soil moisture
676 and runoff in coastal northern california. *Journal of Hydrometeorology*, 14 (2), 443–
677 459.
- 678 Ralph, F. M., Neiman, P. J., Kiladis, G. N., Weickmann, K., & Reynolds, D. W. (2011). A
679 multiscale observational case study of a pacific atmospheric river exhibiting
680 tropical—extratropical connections and a mesoscale frontal wave. *Monthly Weather*
681 *Review* , 139 (4), 1169–1189.
- 682 Ralph, F. M., Neiman, P. J., Wick, G. A., Gutman, S. I., Dettinger, M. D., Cayan, D. R., &
683 White, A. B. (2006). Flooding on california's russian river: Role of atmospheric rivers.
684 *Geophysical Research Letters*, 33 (13).
- 685 Ralph, F.M., Dettinger, M.D., Cairns, M.M., Galarnau, T.J. and Eylander, J. (2018).
686 Defining “atmospheric river”: How the Glossary of Meteorology helped resolve a
687 debate. *Bulletin of the American Meteorological Society*, 99(4), pp.837-839.
- 688 Rutz, J. J., Steenburgh, W. J., & Ralph, F. M. (2014). Climatological characteristics of
689 atmospheric rivers and their inland penetration over the western United states.
690 *Monthly Weather Review* , 142 (2), 905–921.
- 691 Rutz, J.J., Shields, C.A., Lora, J.M., Payne, A.E., Guan, B., Ullrich, P., O'brien, T., Leung,
692 L.R., Ralph, F.M., Wehner, M. and Brands, S. (2019). The atmospheric river tracking
693 method intercomparison project (ARTMIP): Quantifying uncertainties in atmospheric
694 river climatology. *Journal of Geophysical Research: Atmospheres*, 124(24),
695 pp.13777-13802.
- 696 Shields, C. A., Rosenbloom, N., Bates, S., Hannay, C., Hu, A., Payne, A. E., . . Truesdale,
697 J. (2019). Meridional heat transport during atmospheric rivers in high-resolution
698 cesm climate projections. *Geophysical Research Letters*, 46 (24), 14702–14712.
- 699 Shields, C. A., Rutz, J. J., Leung, L. R., Ralph, F. M., Wehner, M., O'Brien, T., & Pierce,
700 R. (2019). Defining uncertainties through comparison of atmospheric river tracking
701 methods. *Bulletin of the American Meteorological Society*, 100 (2), ES93–ES96.



- 702 Solmon, F., Mallet, M., Elguindi, N., Giorgi, F., Zakey, A., & Konaré, A. (2008). Dust
703 aerosol impact on regional precipitation over western Africa, mechanisms and
704 sensitivity to absorption properties. *Geophysical Research Letters*, 35(24).
- 705 Taylor, K. E. (2001). Summarizing multiple aspects of model performance in a single
706 diagram. *Journal of geophysical research: atmospheres*, 106(D7), 7183-7192.
- 707 Ullrich, P. A., & Zarzycki, C. M. (2017). Tempestextremes: A framework for scale-
708 insensitive pointwise feature tracking on unstructured grids. *Geoscientific Model*
709 *Development*, 10 (3), 1069–1090.
- 710 Ullrich, P. A., Zarzycki, C. M., McClenny, E. E., Pinheiro, M. C., Stansfield, A. M., & Reed,
711 K. A. (2021). Tempestextremes v2. 1: A community framework for feature detection,
712 tracking and analysis in large datasets. *Geoscientific model development*
713 *discussions*, 2021 , 1–37.
- 714 Vo et al., (2024). xCDAT: A Python Package for Simple and Robust Analysis of Climate
715 Data. *Journal of Open Source Software*, 9(98), 6426,
716 <https://doi.org/10.21105/joss.06426>
- 717 Wilks, D. S. (2011). *Statistical methods in the atmospheric sciences* (Vol. 100). Academic
718 press.
- 719 Zarzycki, C. M., Ullrich, P. A., & Reed, K. A. (2021). Metrics for evaluating tropical
720 cyclones in climate data. *Journal of Applied Meteorology and Climatology*, 60 (5),
721 643–660.
- 722 Zhang, L., Zhao, Y., Cheng, T. F., & Lu, M. (2024). Future changes in global atmospheric
723 rivers projected by CMIP6 models. *Journal of Geophysical Research: Atmospheres*,
724 129, e2023JD039359
- 725 Zhao, A., Ryder, C. L., & Wilcox, L. J. (2022). How well do the CMIP6 models
726 simulate dust aerosols?. *Atmospheric Chemistry and Physics*, 22(3), 2095-2119.
- 727 Zhu, Y., & Newell, R. E. (1998). A proposed algorithm for moisture fluxes from
728 atmospheric rivers. *Monthly weather review* , 126 (3), 725–735

729
730
731

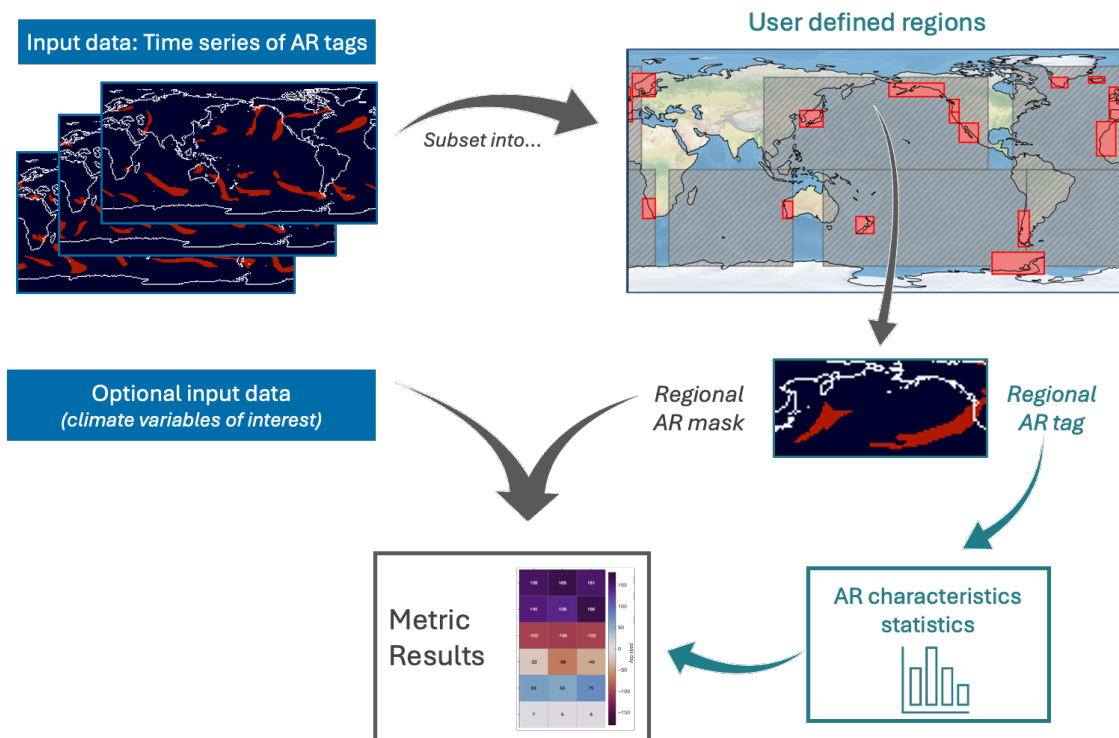
732 **Supplementary information**

733 In a separate document

734
735

736 **Figures and tables**

737



738

739

740 Fig. 1. AR metric tool workflow. Input data include time slices of AR
741 tags from ARDTs of user choice, and optional climate data
742 associated with ARs. The data are then subset into user-defined
743 rectangular domains (blue boxes for ocean basins, red boxes for
744 landfall regions) for regional tags and masks. User preferred
745 statistical tools are applied on the regional AR tags to obtain AR
746 characteristics. Finally, AR characteristics and AR masked climate
747 data are presented as metric results.

748

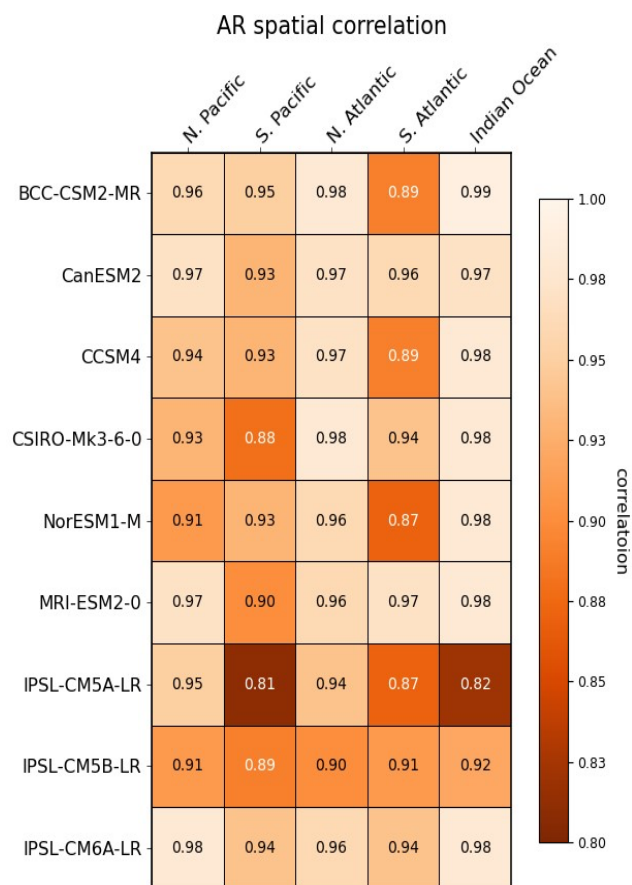


749
 750
 751
 752
 753
 754
 755

Table 1. List of AR metrics and diagnostics in this study. Numbers in the table indicate the number of regions where the metrics are applied. Each column is one AR property. Underscored items are model evaluation metrics, items in italic form are diagnostics of AR properties.

metrics/ diagnostics	ARs over Ocean Basins						Landfalling ARs		
	frequency	central latitude	central longitude	size	width	length	counts (frequency)	peak day	precipitation
<u>mean bias</u>	5	5	5	5	5	5	16	16	16
<u>spatial correlation</u>	5								
<u>IoU</u>							16		
<i>spatial distribution</i>	5						16		
<i>sampling histogram</i>		5	5	5	5	5			
<i>monthly climatology histogram</i>							16		

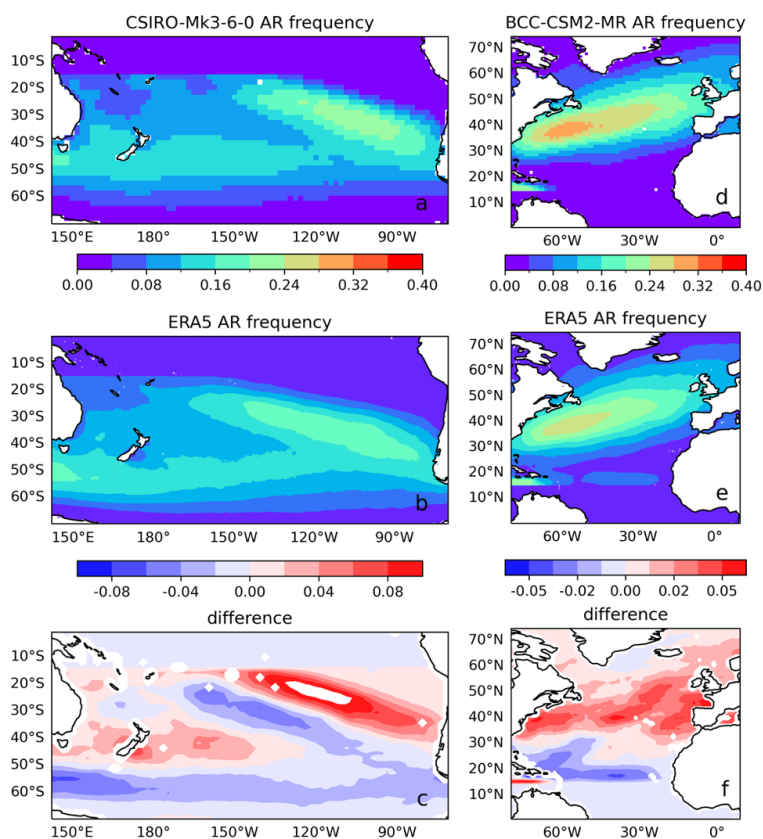
756



757

758 Fig. 2. Spatial pattern correlation of AR frequency for the period 1979-
 759 1989 between ERA5 and climate models for major ocean basins.

760



761
762
763
764
765
766
767
768
769

Fig. 3. AR frequency in the South Pacific for (a) CSIRO-MK3-6-0, (b) ERA5 and their difference (c) as (a) - (b). AR frequency in the North Atlantic for (d) BCCSM2-MR, (e) ERA5 and their difference (f) as (d) - (e)

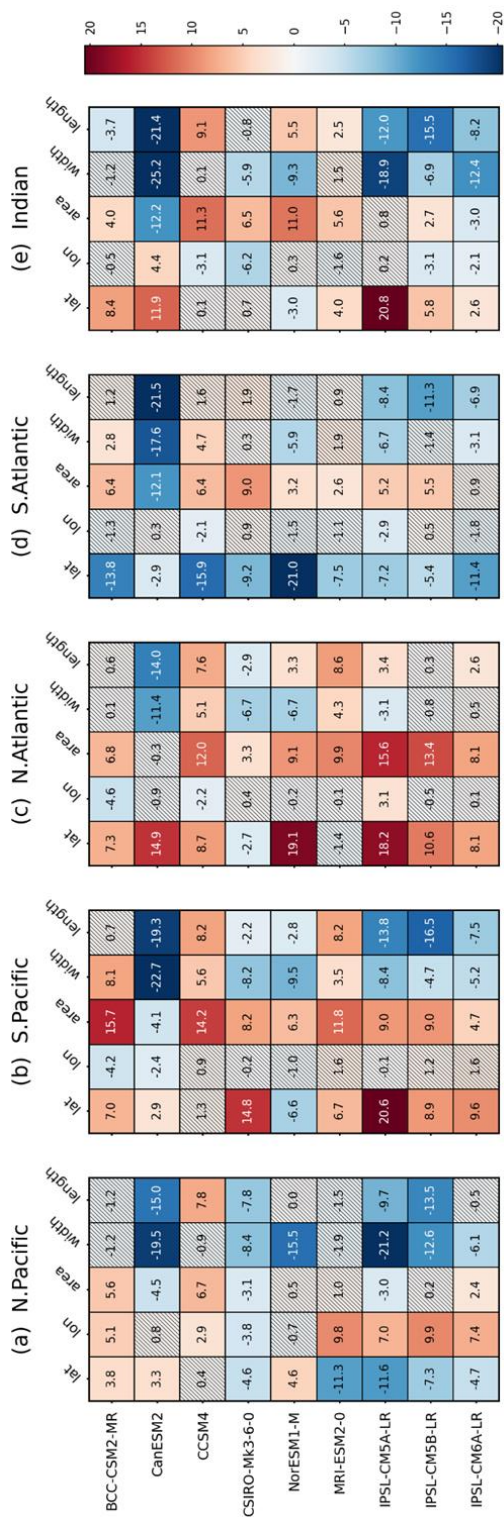
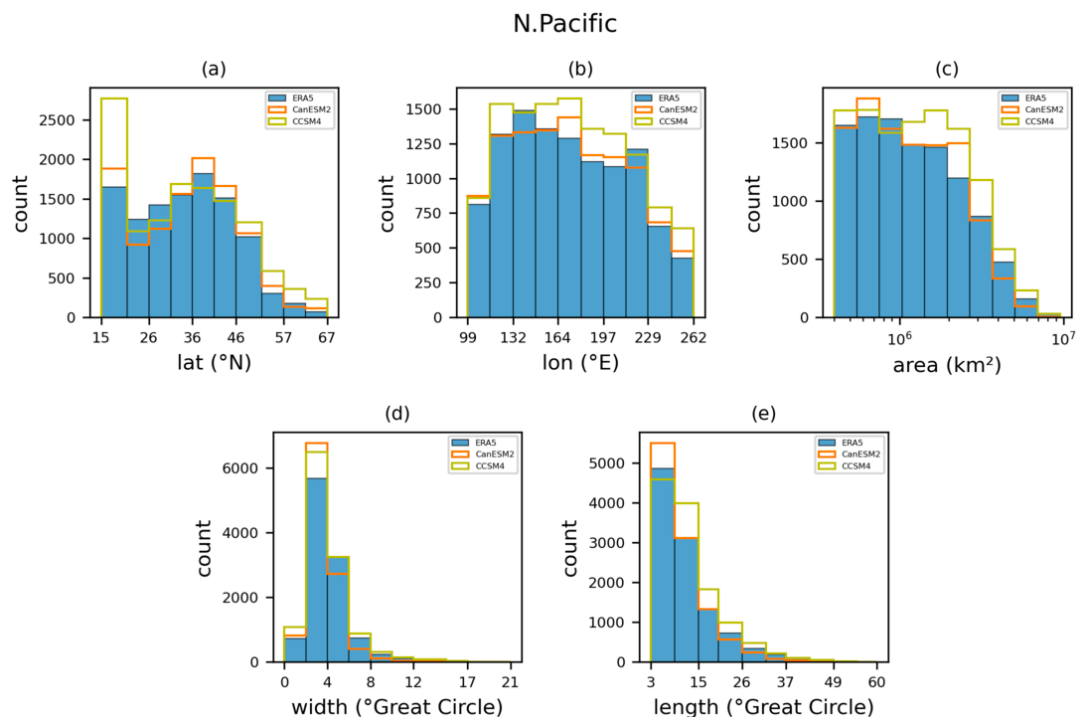


Fig. 4 AR characteristics bias (normalized as Z-score) in climate models for major ocean basins. Hatching indicates that the differences are statistically insignificant.



771

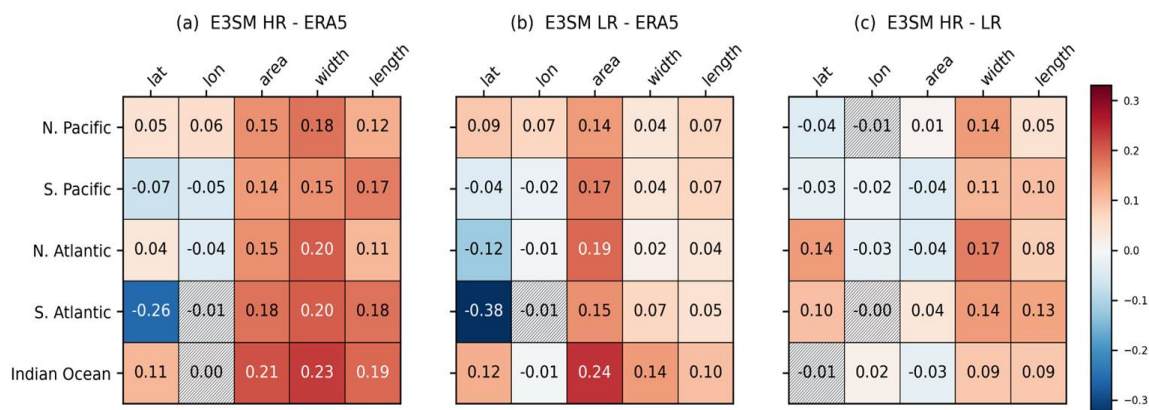


772

773 Fig. 5 North Pacific AR characteristics distribution for (a) central latitude, (b)
774 central longitude, (c) area, (d) width and (e) length, in ERA5 reanalysis,
775 CanESM2 and CCSM4 model



776
 777



778
 779
 780
 781
 782
 783
 784
 785

Fig. 6. AR characteristics bias in E3SM (a) HR and (b) LR simulations. (c) is the difference between HR and LR. Hatching indicate that the differences are statistically insignificant.

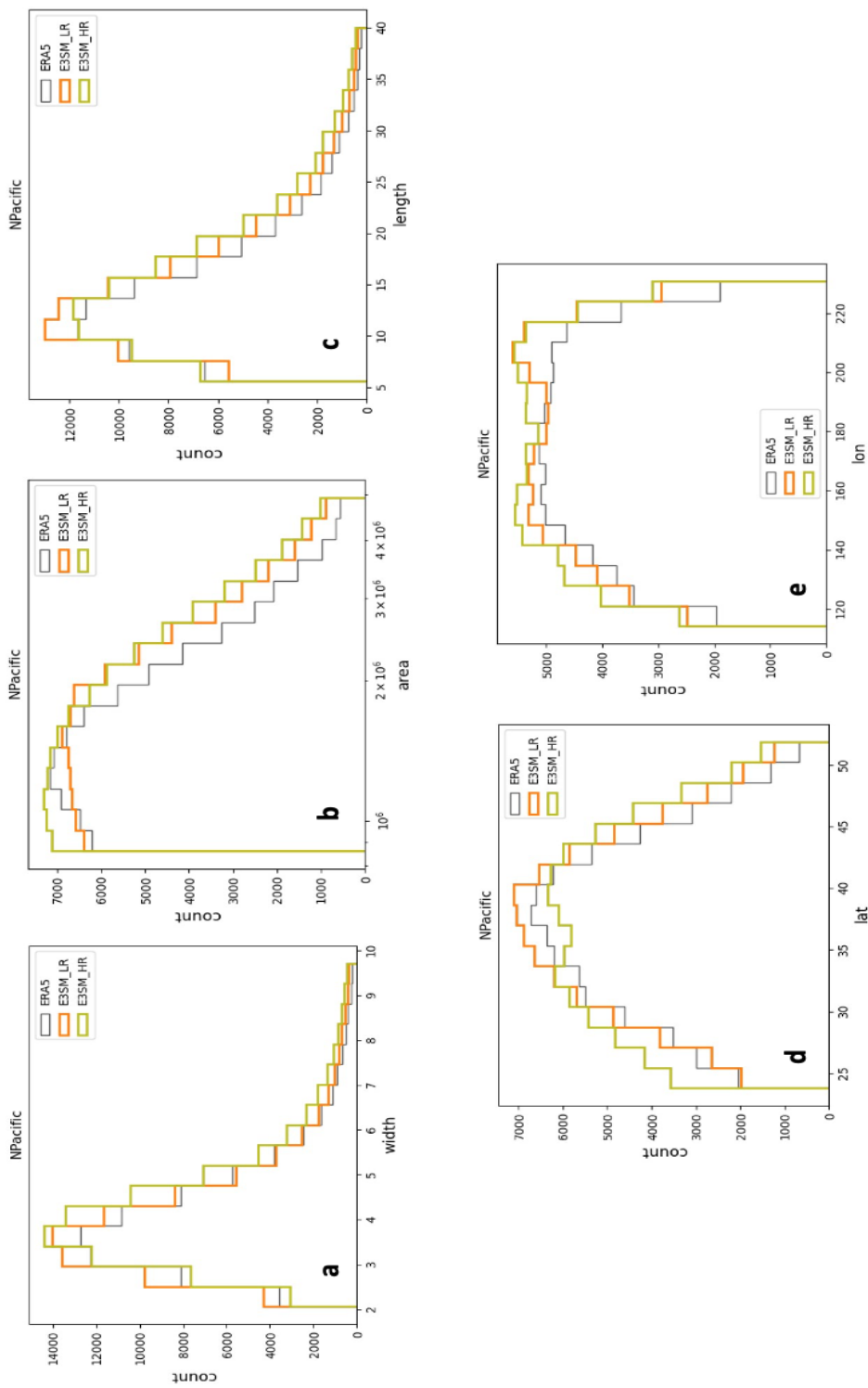
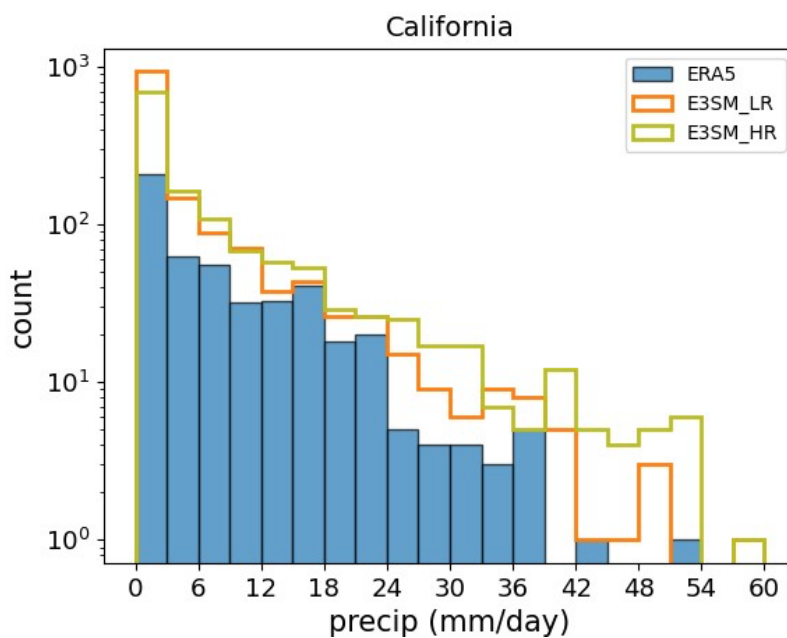


Fig. 7. AR characteristics distribution of (a) width (° great circle), (b) area (km²), (c) length (° great circle), (d) central latitude (°N), and (e) central longitude (°E) in the North Pacific for ERA5, E3SM LR and LR simulations.



787
788
789
790
791



792
793
794
795
796
797
798

Fig. 8. Landfalling AR precipitation histogram in California from 1990-1999 in the ERA5 reanalysis, E3SM HR and LR simulations.



799
800
801
802
803
804
805
806
807

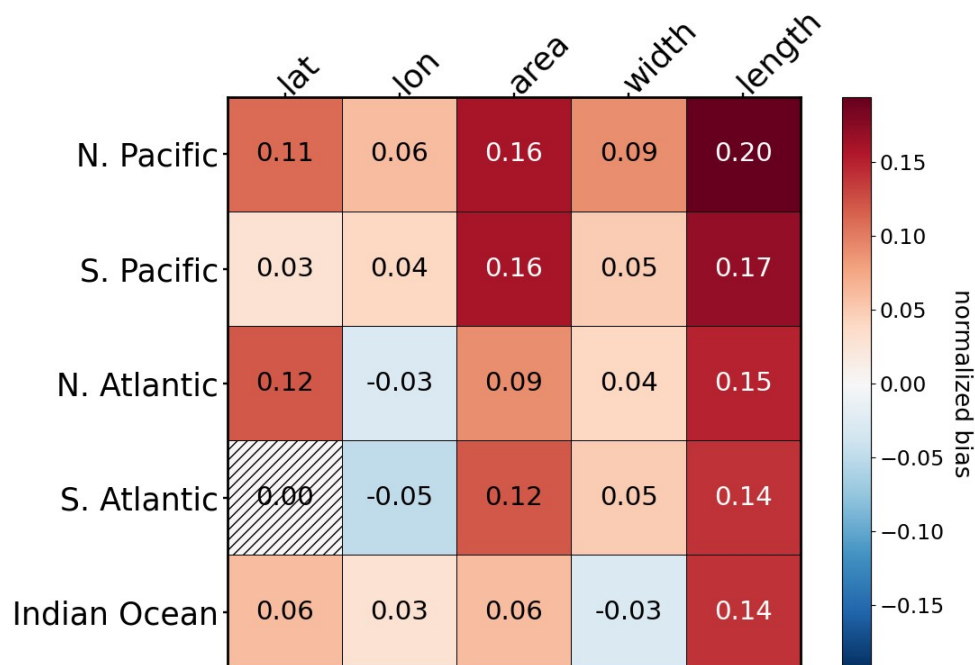
Table 2. AR landfall concurrence in Mundhenk and TE, normalized by total counts of AR landfalls detected in both ARDTs for different regions. Values are shown in percentage.

Region	California	S. America	N. Europe	Australia	S. Africa	Baja	Pacific Northwest	New Zealand
Concurrence (%)	56	68	82	62	51	30	72	77
Region	Alaska	UK	W. Europe	Iceland	Greenland	E. Asia	Antarctica	New England
Concurrence (%)	81	84	74	77	72	56	69	83

808



809
810
811

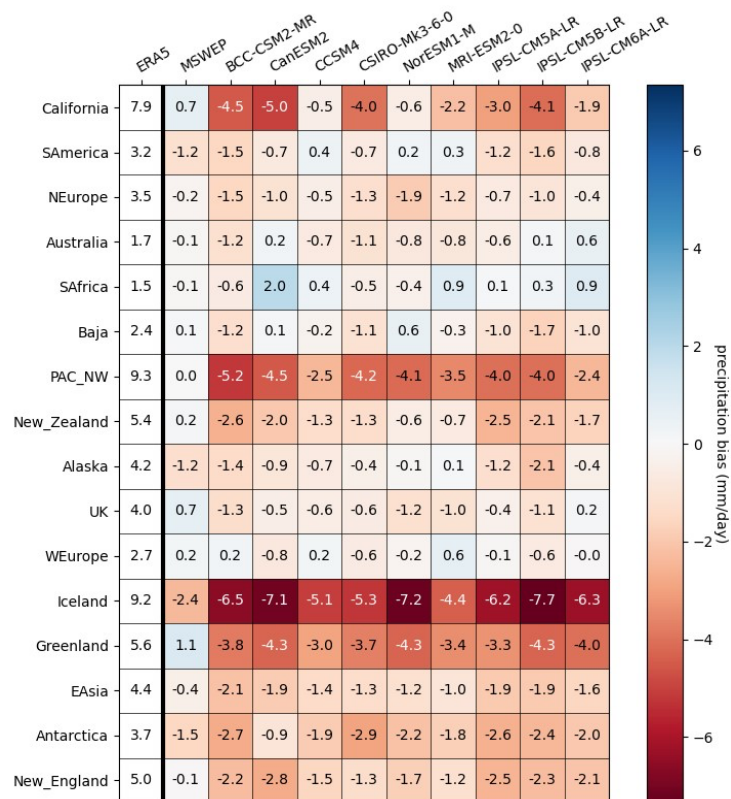


812
813
814
815
816
817
818
819

Fig. 9. AR characteristic difference between Mundhenk and TE in ERA5



820
 821

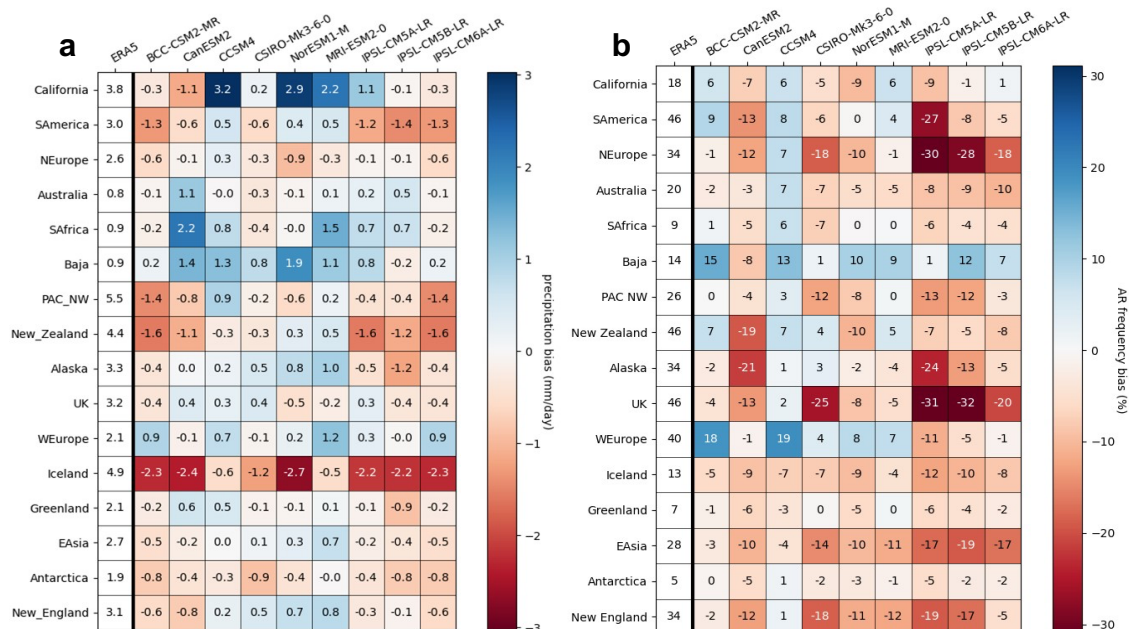


822
 823
 824
 825
 826
 827

Fig. 10. Landfalling AR precipitation bias in climate models relative to ERA5 (the first column). The MSWEP data is also included in the second column as an additional reference data, showed as the difference between ERA5.



828
 829

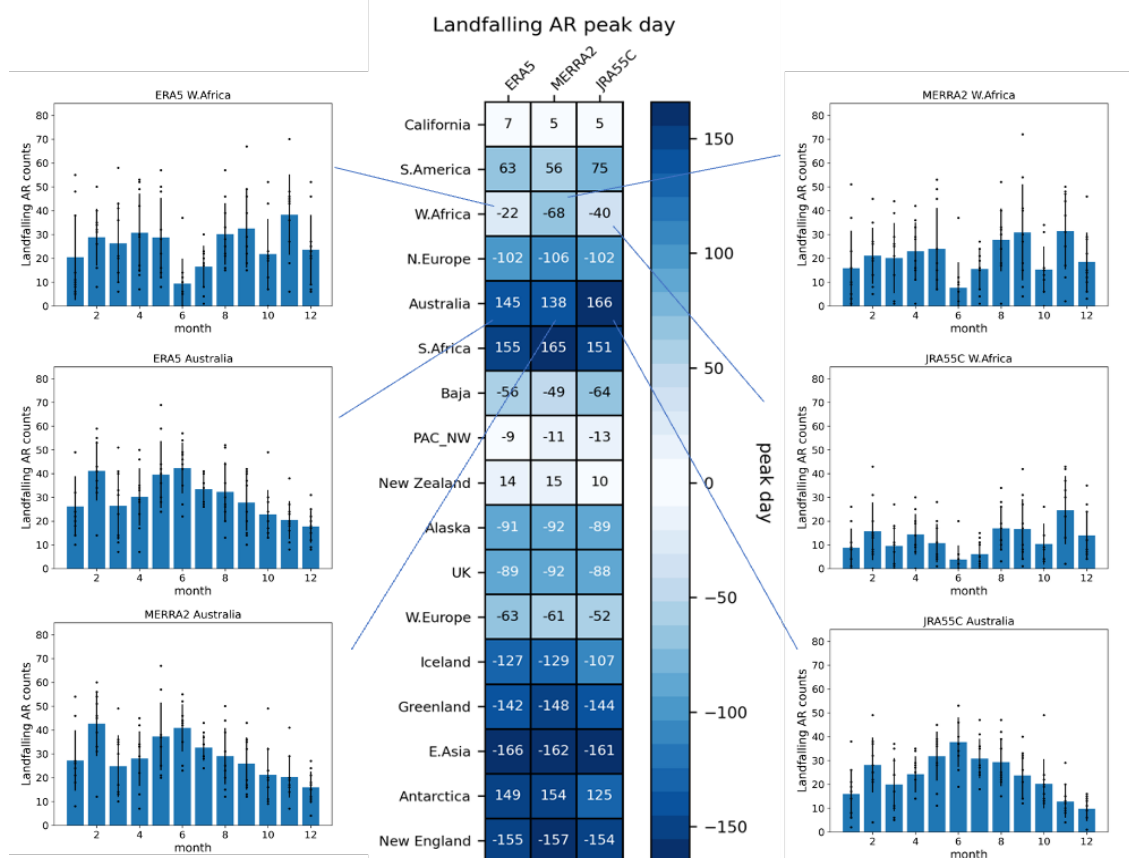


830
 831
 832
 833
 834

Fig. 11. (a) Total precipitation bias and (b) landfalling AR frequency bias



835



836

837

838

839

840

841

842

843

844

845

846

847

848

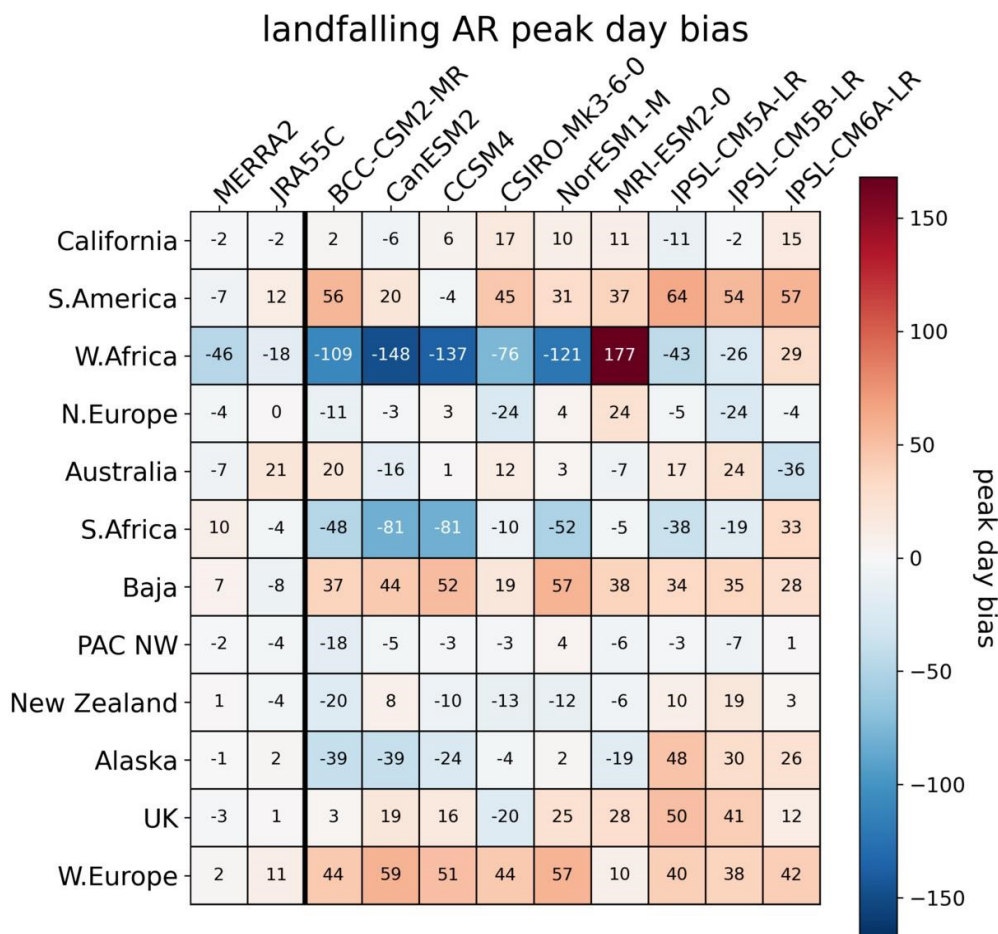
849

850

Fig. 12 (a) Landfaling AR peak day in ERA5, MERRA2, and JRA55C reanalysis. (b-g) show examples of probability distribution. Height of the blue bars indicate the time mean counts. Black dots represent peak day for each individual year, and vertical bars are the standard deviation range in the 10-year data from 1979-1988

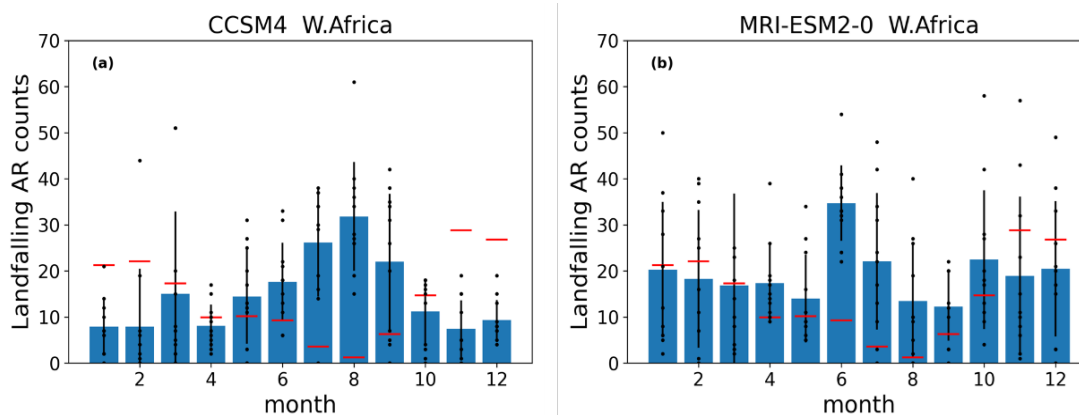


851
 852



853
 854
 855
 856
 857
 858
 859
 860
 861
 862

Fig. 13. Landfalling AR peak day bias in reanalyses and models compared with ERA5.



863
864
865
866
867
868
869
870
871

Fig. 14. Landfalling AR counts in (a) CCSM4 and (b) MRI-ESM2-0 for western Africa region. Height of the blue bars indicate the time mean counts. Vertical lines represent the standard deviations. Black dots represent counts for each individual year. Red bars show ERA5 values as the reference.