

1 **A new metrics framework for quantifying and intercomparing atmospheric rivers**
2 **in observations, reanalyses and climate models**

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

- 6
7 1. Lawrence Livermore National Laboratory, Livermore, CA, USA
8 2. Department of Earth and Atmospheric Sciences, Indiana University,
9 Bloomington, IN, USA
10 3. Climate and Ecosystem Sciences Division, Lawrence Berkeley National Lab,
11 Berkeley, CA, USA

12
13
14 Correspondence to: Bo Dong (dong12@llnl.gov)

15
16
17 **Key points:**

- 18 1. A metrics package designed for easy analysis of AR characteristics and statistics is
19 presented
20 2. The tool is efficient for diagnosing systematic AR bias in climate models, and useful
21 for evaluating new AR characteristics in model simulations
22 3. In climate models, landfalling AR precipitation shows dry biases globally, and AR
23 tracks are farther poleward (equatorward) in the north and south Atlantic (south Pacific
24 and Indian Ocean)

25
26
27 **Abstract**

28
29 We present a new atmospheric river (AR) analysis and benchmarking tool, namely
30 Atmospheric River Metrics Package (ARMP). It includes a suite of new AR metrics that
31 are designed for quick analysis of AR characteristics via statistics in gridded climate
32 datasets such as model output and reanalysis. This package can be used for climate
33 model evaluation comparing with reanalysis and observational products. Integrated
34 metrics such as mean bias and spatial pattern correlation are efficient for diagnosing
35 systematic AR biases in climate models. For example, the package identifies that in
36 CMIP5 and CMIP6 models, AR tracks in the south Atlantic are positioned farther
37 poleward compared to ERA5 reanalysis, while in the south Pacific, tracks are generally
38 biased towards the equator. For the landfalling AR peak season, we find that most
39 climate models simulate a completely opposite seasonal cycle over western Africa. This
40 tool can also be used for identifying and characterizing structural differences among

41 different AR detectors (ARDTs). For example, ARs detected with the Mundhenk
42 algorithm exhibit systematically larger size, width and length compared to the
43 TempestExtremes (TE) method. The AR metrics developed from this work can be
44 routinely applied for model benchmarking and during the development cycle to trace
45 performance evolution across model versions or generations and set objective targets
46 for the improvement of models. They can also be used by operational centers to
47 perform near real-time climate and extreme events impact assessment as part of their
48 forecast cycle.

49
50

51 **1. Introduction**

52

53 Atmospheric rivers (ARs) are dynamically driven, synoptic-scale filamentary structures
54 of water vapor jets that play important roles in the global water cycle and regional
55 weather and hydrology (Ralph et al. 2013; Gimeno et al. 2014; Shields et al. 2019;
56 Payne et al. 2020; O'Brien et al., 2022). These narrow, concentrated corridors of
57 moisture in the atmosphere can carry an immense amount of water, often compared to
58 the flow of multiple major rivers combined (Ralph and Dettinger, 2011), and account for
59 more than 90% of extratropical poleward water vapor transport (Zhu and Newell, 1998;
60 Newman et al. 2012; Ullrich et al. 2021). When making landfall or interacting with
61 topography, ARs can produce extreme weather, including heavy rainfall and strong
62 winds, in turn leading to severe flooding and landslides. These effects can devastate
63 natural landscapes, agricultural fields, human infrastructure, and disrupt businesses and
64 services, leading to significant economic losses (Ralph et al., 2006; Leung and Qian,
65 2009; Neiman et al., 2011; Neiman et al., 2013; Gershunov et al., 2017). However, ARs
66 are also essential for delivering water for agriculture, ecosystems and human
67 consumption; in the western United States alone they are responsible for one-third to
68 one-half of total annual precipitation (Ralph and Dettinger, 2011).

69

70 Because ARs can be responsible for both beneficial and detrimental impacts,
71 understanding and modeling of these features, particularly in light of climate change, is
72 an important topic. To date, however, proposed definitions of ARs have yet to be
73 widely adopted (Ralph et al., 2018), which has in turn made it difficult to draw
74 conclusions about how these features may be changing. Numerical algorithms for
75 objective identification of ARs, namely AR detectors (ARDTs) (e.g., Neiman et al., 2009;
76 Dettinger, 2011; Ralph et al., 2013; Mundhenk et al. 2016; Ullrich and Zarzycki 2017;
77 Ullrich et al., 2021) have widely facilitated broader studies of AR characteristics and
78 impacts (Shields et al., 2018; Rutz et al., 2019; O'Brien et al., 2022). However, as
79 ARDTs are usually designed with particular research questions in mind, the lack of a
80 unified framework that is applicable to different ARDTs in a collective way has

81 challenged the benchmarking and intercomparison of the models' representation of
82 ARs. The analysis workflow and code in one study cannot be easily applied in another
83 study using a different ARDT. Consequently, studies like intercomparison of ARDTs, or
84 analysis based on an ensemble of ARDTs cannot be readily executed without extensive
85 collaboration or community efforts. In addition, research of this kind cannot be easily
86 repeated or updated when newer versions of ARDTs have been developed, or newer
87 observational data products have become available. As such, a universal analysis
88 framework that is independent of ARDT is in demand in our AR research community.

89
90 Within AR research, one major branch focuses on evaluating the performance of
91 forecast or climate models in simulating ARs. Since the number of climate models under
92 active development and used in the research community has increased substantially in
93 recent decades, with many supporting multiple configurations and parameterization
94 choices, routine evaluation of ARs during model development lifecycles requires a
95 quantitative climate data assessment evaluation workflow that allows comparing AR
96 characteristics from different ARDTs. We believe progress in improving our
97 understanding of ARs and their impacts could be accelerated with a dedicated tool for
98 calculating AR statistics and evaluation metrics in climate models and gridded data
99 products. Preferably, such an analysis tool should be seamlessly applicable to multiple
100 data sources (including observations, forecast, reanalysis and different models) with
101 simply a few commands, minimizing users' efforts to manage inconsistencies in data
102 format, coordinate system and spatial coverage of different data.

103
104 In this paper, we propose a new AR analysis framework that includes a diverse suite of
105 metrics that is designed for easy quantification of AR characteristics and statistics in all
106 types of gridded climate data, with the expectation that such a metric suite would be
107 efficient for ARDT intercomparison and climate model evaluation. Following the
108 introduction, section 2 describes the general design and workflow of the AR metrics
109 tool. Section 3 presents several model evaluation and ARDT intercomparison
110 application examples using the metrics evaluation package. Discussion and future
111 development plans are in section 4.

112 113 **2. AR metrics package design and workflow**

114 115 **2.1 Metrics workflow**

116
117 Figure 1 shows the general design and workflow of the AR Metrics Package (ARMP).
118 The input data includes AR objects and optional climate variables of relevance to ARs,
119 such as precipitation, winds, and temperature. The AR tags can be produced by any
120 regional or global ARDT, including those based on relative (e.g., TempestExtremes or

121 TE; Ullrich and Zarzycki 2017; Ullrich et al. 2021), fixed-relative (e.g., Mundhenk_v3;
122 Mundhenk et al. 2016), and absolute (e.g., Lora_v2; Lora et al. 2017) thresholds on the
123 moisture field.

124
125 AR metrics are calculated in user-defined geographic domains. The upper right panel in
126 Fig. 1 shows examples of regions that were selected for landfalling AR diagnostics (red
127 boxes in the panel, lat-lon boundaries are listed in the supplementary table S3). These
128 regions, mostly located along the west coast of continents, are known to have frequently
129 observed AR landfalls (Guan and Waliser 2015, Algarra et al. 2020). We purposely use
130 rectangular region boundaries for simplicity and to avoid masking files; numerous tools
131 are already available for sub-selection of data using latitude-longitude boundaries.

132
133 Apart from metrics for AR landfall regions, rectangular region subsetting is also useful
134 for analyzing AR geometric features over global oceans. Currently there are 5 ocean
135 basins pre-defined in the framework – the North Pacific, South Pacific, North Atlantic,
136 South Atlantic, and South Indian Ocean (blue boxes in Fig. 1 upper right panel; lat-lon
137 coordinates in table S3 in the supplement). The analysis domains are fully customizable
138 with specified latitude and longitude boundaries from local (depending on spatial data
139 resolution) to global scales.

140
141 The regional segments of AR tags can then be used standalone as regionally cropped
142 AR objects for feature metrics calculation. In this paper we provide five metrics
143 application examples in section 3, such as AR geometry, frequency and landfall
144 seasonality. Alternatively, the regional tags can be used as masks for AR associated
145 weather, dynamical and thermodynamical processes. An example of evaluating AR
146 precipitation in climate models is given section 3.

147
148 For AR geometrical metrics, statistics gauging the consistency of latitude, longitude,
149 width, length, and size are required as intermediate input data to the workflow. In the
150 examples presented in this paper, we use the 'BlobStats' tool (Ullrich et al. 2021) to
151 calculate these statistics, where latitude and longitude are weighted by the moisture
152 field, width and length are based on principle component analysis (PCA; Inda-Díaz et al.
153 2021), and size is based on a count of the number of contiguous grid cells in the
154 feature. This tool can be called and run within the AR metrics framework, albeit
155 requiring an additional installation. Users can also optionally use their preferred
156 statistical tool for AR geometry calculation and then feed the data back to the metrics
157 workflow.

158
159 The metrics and diagnostics are integrated into the framework, which can be
160 customized and expanded subject to the objective of research. Table 1 lists all the AR

161 metrics and diagnostics used in this study. The AR metrics are composed of AR
162 properties (as shown in the top row) and evaluation metrics. Similarly, the AR
163 diagnostics are composed of AR properties and statistical diagnostics. The number of
164 regions that these metrics are applied to are indicated by the numbers in the table.
165

166

167 **2.2 Software structure, coding environment and data format**

168

169 The metrics code is open-source and python-based, and it handles gridded AR tag and
170 climate data using Xarray (<https://xarray.pydata.org>, Hoyer et al., 2017) and its
171 extension xCDAT package (Xarray Climate Data Analysis Tools,
172 <https://xcdat.readthedocs.io>, Vo et al., 2024). It also leverages several utility functions in
173 the PCMDI Metrics Package (PMP; Lee et al. 2024), such as the regional regridding
174 tool, land-sea mask, and portrait plot. These packages are compatible with one another,
175 readily available and easy to install. The code repository can be accessed at
176 <https://github.com/PCMDI/ARMP>, and relevant wiki documents including a demo Jupiter
177 notebook are provided with installation instructions and application examples.
178

178

179 The code consists of 7 major components: workflow controller, I/O, data QAQC,
180 functional utilities, regional statistics, benchmarking metrics and graphics. It accepts AR
181 masks and climate data files in NetCDF format as input data. Input filenames are listed
182 in a pointer file as a configuration parameter to the metrics package. Output files are in
183 NetCDF format for intermediate and diagnostic outputs, and JSON format for computed
184 metrics. The regional statistics module integrates a few commonly used statistics for AR
185 properties (e.g., AR frequency and AR precipitation), and newly developed statistics
186 (e.g., AR landfall peak day). External statistical tools, e.g., BlobStats for calculating AR
187 geometry, can also be called from this package. These statistics are then fed into the
188 metrics module. AR metrics included in this framework are described in section 2.3.
189

189

190 The metrics tool can be applied to data with different resolution, domain (e.g., a list of
191 data files with mixed global and regional spatial extent), and coordinate system (e.g.,
192 180° or 360° longitude coordinates; monotonically decreasing latitude coordinates),
193 minimizing the effort required to prepare the input data files. It is compatible with CF-
194 compliant NetCDF files as well as some non-compliant data structures. It also aims to
195 intelligently flag imperfect data, including files with corrupted data values, or with
196 incorrect datetime calendar.
197

197

198

199 **2.3 AR benchmarking metrics**

200

201 Metrics have been widely used to quantify climate model performance (Taylor 2001;
202 Gleckler et al. 2008; Wilks 2011; Zarzycki et al. 2021). In the AR community, a set of
203 common metrics are also increasingly employed over the past few years, such as mean
204 bias (Guan and Waliser 2017; Chapman et al. 2019), weighted ensemble mean bias
205 (Massoud et al. 2019), RMS error and relative RMS error (Guan and Waliser 2017),
206 spatial pattern correlation (Chapman et al. 2019; Huang et al. 2021), ratio of spatial
207 standard deviation (O'Brien et al. 2022), and skill scores for assessing AR predictions
208 (Wick et al. 2013, Nardi et al. 2018) and model performance (Zhang et al. 2024). While
209 these quantitative measures are case-specific and depend on the aim of these studies,
210 there is value in synthesizing commonly used metrics into one comprehensive analysis
211 tool. Here we describe a suite of diverse metrics used in this study, including both
212 commonly-used as well as newly proposed metrics.

213

214 2.3.1 Mean bias

215

216 We use mean bias to measure how close a climate data product is to an appropriately
217 chosen reference dataset. The statistical significance of the mean bias is measured
218 using the Z-test. For sake of completeness, the mathematical formula of mean bias and
219 z-score are given in the Appendix. Under this test, the difference between the means of
220 two samples is considered to be statistically significant at the 95% confidence level if
221 the magnitude of the z-score is greater than 1.96. When comparing across different
222 variables, a commonly-used measure is the normalized bias, with the data normalized
223 by the standard deviation of the reference field. In this study, we simply use z-score as
224 the normalized bias, as it incorporates both bias and statistical significance in one
225 succinct formula.

226

227 2.3.2 Spatial pattern similarity

228

229 The spatial pattern correlation is a measure used to quantify the similarity between two
230 spatial fields without reflecting the magnitude of the difference. Here we compute the
231 spatial pattern correlation using the Pearson correlation coefficient. The statistical
232 significance of correlation is determined by the two-tailed p-value of the cumulative
233 distribution function (CDF) of the t-statistic. The mathematical formula for the Pearson
234 correlation coefficient and its corresponding significance test is given in the Appendix.
235 Given that ARs have notable seasonal and interannual latitudinal shifts, we propose a
236 new method to estimate the effective sample size n_e as the number of Principal
237 Component Analysis (PCA) modes required to explain more than 95% of the total
238 variance in the AR tag data. The cumulative variance explained by the principal
239 components is expressed as

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

240
241

242 where the λ_i are the eigenvalues of the spatial correlation matrix of the data, and p is the
243 total number of principal components. Estimating n_e based on ERA5 reanalysis data, we
244 find that the effective sample sizes for spatial pattern correlation are generally small,
245 ranging from 14 - 27 PCs necessary to explain more than 95% of total variance for the 5
246 ocean basins (Table S4 in supplementary information).

247

248 2.3.3 Temporal detection similarity

249

250 The AR binary occurrence time series is a time series variable equal to one when an AR
251 is present in a given region and zero otherwise. The overlap between two AR
252 occurrence time series is measured by the Intersection over Union (IoU) metric. The
253 metric is written as

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

254

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

257

258

259 3. Metrics applications

260

261 In this section, we present five example applications using the metrics tool for assessing
262 ARs in climate models, including evaluation of AR frequency and characteristics,
263 comparison of ARs in high- and low-resolution simulations, sensitivity of ARs to choice
264 of ARDT, precipitation bias associated with ARs and landfalling AR seasonality.

265

266 3.1 AR tag and climate data

267

268 We compare the TE ARDT on the 6-hourly integrated water vapor transport (IVT) data
269 from three reanalysis products - ERA5 (Hersbach et al. 2020), MERRA-2 (Gelaro et al.
270 2017) and JRA-55C (Japan Meteorological Agency, Japan 2015) to obtain AR tags for
271 reanalyses. Given its longer data record and finer model resolution, we use ERA5 as
272 the default reference in this study. To demonstrate how results are sensitive to the
273 choice of ARDTs, we then use the fixed-relative (Mundhenk_v3) tags from ERA5 data.

274

275 To evaluate ARs in climate models, we use the archived AR tags from the Atmospheric
276 River Tracking Method Intercomparison Project (ARTMIP) Tier 2 experiment, which is
277 based on the coupled CMIP model simulations for the historical and 21st century
278 projection periods (Shield et al. 2019, Rutz et al 2019, O'Brien et al, 2022). The tag data
279 include six of the CMIP5 models (CCSM4, CSIRO-Mk3-6, CanESM2, IPSL-CM5A-LR,
280 IPSL- CM5B-L, and NorESM1-M) and 3 of the CMIP6 models (BCC-CSM2-MR, IPSL-
281 CM6A-LR, MRI-ESM2-0). Grid information of these models is listed in supplementary
282 Table S1. All the tags data are in 6-hourly temporal frequency. For model evaluation
283 purposes in our application examples, only TE tags from the archive are selected.
284

285 We further use simulations from the Energy Exascale Earth System Model (E3SM;
286 Golaz et al. 2019, Caldwell et al. 2019) high resolution (HR, 0.25°, ~28 km grid) and low
287 resolution (LR, 1°, ~111 km grid) experiments to examine the sensitivity of ARs to
288 model resolution. Comparison of the grid parameters of the two models is also shown in
289 supplementary Table S2. Except for their different horizontal grid spacing, both E3SM-
290 HR and E3SM-LR use an identical set of physical parameters, and the simulations
291 follow a similar protocol of the Coupled Model Intercomparison Project Phase 6 (CMIP6;
292 Eyring et al. 2016).
293
294

295 **3.2 Basic AR characteristics in CMIP5 and CMIP6 models**

296 297 3.2.1 AR frequency

298
299 We first analyze the pattern of AR occurrence frequency over a 10-year period (1979-
300 1988) for the five major ocean basins from section 2.2. From the spatial distribution of
301 the AR frequency, we calculate the pattern correlation between selected climate models
302 and ERA5. The spatial pattern correlation coefficient is shown in Fig. 2. Notably the
303 correlations are statistically significant for all models and regions. This suggests that
304 climatologically, all climate models simulate AR density and spatial distribution that
305 broadly resemble reanalysis on planetary scale. This is evidenced in the spatial AR
306 occurrence density maps in Fig. 3 (a-b) and (d-e).
307

308 The high spatial correlation (e.g., in Fig. 3, $r = 0.88$ in S. Pacific and $r = 0.98$ in N.
309 Atlantic) is mainly a result of the similar spatial gradient (as in Fig 3a-b, and Fig 3d-e) of
310 the AR frequencies, rather than a similar magnitude of frequency at each grid point in
311 the two datasets. For instance, if the AR frequency values in one map are doubled
312 compared to those on the other map, the spatial patterns, or spatial structures of the
313 two, can still be perfectly correlated. Since climatologically ARs are largely clustered
314 along the storm track, with nearly no occurrence over a large portion of the basin

315 domain, it is natural that the pattern correlations are significant in most cases. Similar
316 high pattern correlations of AR frequencies are also noted in other studies (e.g., Huang
317 et al. 2020; Guan et al. 2023). In other words, the spatial correlation coefficient is not
318 that indicative for the magnitude resemblance of the AR spatial frequency. Therefore,
319 these metric results can be better interpreted together with AR frequency maps with
320 spatial gradient.

321
322 While the spatial correlation coefficient synthesizes the level of pattern consistency,
323 difference maps further reveal spatial discrepancies. For example, Fig. 3c shows that
324 South Pacific AR tracks shift farther towards the equator in the CSIRO model than in
325 ERA5. While in the North Atlantic basin (Fig. 3f), AR tracks are displaced more
326 poleward in the BCC model. The further north AR location is likely associated with the
327 poleward jet stream bias in CMIP6 models (Bracegirdle et al. 2020; Harvey et al. 2020).
328 Another example is the AR frequency distribution over the Indian Ocean for BCC-CSM-
329 MR (Fig. 3g-i) and IPSL-CM5A-LR (Fig. 3j-l) model. Even though, compared to ERA5,
330 both models show significant spatial correlation in Fig. 2 ($r=0.99$ and $r=0.82$
331 respectively), the spatial bias pattern in IPSL-CM5A-LR exhibits a more apparent
332 latitudinal shift than in BCC-CSM-MR.

333
334

335 3.2.2 AR geometric features in major ocean basins

336
337 The portrait plots in Fig. 4 show normalized biases (as z-score) of AR characteristics in
338 climate models for the 5 major ocean basins. Several striking results emerge. For
339 instance, in the North Pacific, the CMIP5 and CMIP6 AR geometry, in terms of width
340 and length, are significantly smaller than the ERA5 reanalysis. One possible cause of
341 such biases is that the AR blobs detected with TE in the relatively lower resolution
342 climate models are geometrically less curvy, and less pointy at the ends; for example,
343 supplementary Fig. S2 shows an example time slice of AR blobs in the ERA5 and BCC
344 model. It is clear that the highlighted AR blob in the BCC model exhibits a “cut-off”
345 feature at both ends, thus shorter in length than the ERA5 reanalysis. And although
346 visually the blob is wider, the PCA based width is actually narrower due to its less curvy
347 blob geometry. In contrast, for all other ocean basins, the AR sizes (area) are generally
348 bigger in climate models. The figures also show notable latitudinal model AR biases,
349 such that compared to the reanalysis, ARs tend to shift towards higher latitudes in the
350 North and South Atlantic and biased towards the equator in the South Pacific and Indian
351 Ocean. To assist in understanding of these geographical biases, a set of AR frequency
352 maps over global ocean basins for each climate model are provided in the
353 supplementary material.

354

355 Fig. 4 also helps identify outliers of a specific model or variable. For example, although
356 most climate models tend to simulate larger ARs than observed (indicated by the
357 positive values in the area columns), one notable exception is the CanESM2
358 model [which has significantly smaller AR width, length, and area than other models and
359 ERA5 reanalysis. Taking a closer look into the AR width and length in the North Pacific
360 in Fig. 5, we see that CanESM2 simulates more smaller ($<1.8 \times 10^6 \text{km}^2$) ARs and fewer
361 larger ($>1.8 \times 10^6 \text{km}^2$) ARs than the reanalysis, resulting in negative mean biases. This
362 type of histogram helps us better understand the AR distribution discrepancies.

363
364 Another example is from the CCSM4 model simulations. The higher bounds of the
365 model histogram in nearly all fields indicate that the CCSM4 model simulates more ARs
366 than the reanalysis, with bigger size indicated as taller area bars in Fig. 5c. The higher
367 ARs counts (~ 500 more counts than ERA5) in the model are mostly located in the high
368 latitudes and the tropics south of 20°N (Fig. 5a), spreading across all longitude (Fig. 5b).
369 Fig. 5d and 5e show that the additional ARs in CCSM4 are narrower and/or longer in
370 shape. These differences may arise from various characteristics of the models, such as
371 the dynamical core (e.g., finite volume in CCSM4, T63 triangular spectral truncation in
372 CanESM2, spectral-transform in ERA5), grid resolution (see supplementary Table S1),
373 and the effect of data assimilation (Buizza et al. 2018) in the ERA5 system.

374

375 **3.3 ARs in high and low resolution E3SM simulations**

376

377 We now apply the metrics and diagnostics identified in section 2.3 to E3SM HR and LR
378 simulations. ARs in both HR and LR exhibit similar structural differences compared to
379 the ERA5 (Fig. 6a, b). They are bigger in terms of area, width, and length, and biased
380 towards higher latitudes in the North Pacific and South Atlantic, as indicated by the
381 positive numbers. Zonally, ARs in E3SM are more westward distributed in the North
382 Pacific (positive biases), and more eastward distributed in the North Atlantic and South
383 Pacific (negative biases). One difference we see between the two experiments is that in
384 the North Atlantic basin, AR tracks in the HR are shifted more northward than in the LR
385 simulation.

386

387 Figure 6c shows AR differences between E3SM HR and LR models. The most
388 noticeable differences are that the HR simulates wider and longer ARs than the LR
389 model over all ocean basins. The AR size, in the area column, however, shows mixed
390 results which are not consistent with systematic biases in width and length. This is
391 probably because of different AR geometric properties in the HR and LR simulations.
392 For example, in Supplementary Figure S4, the highlighted AR blob in the North Atlantic
393 is longer but smaller in the LR compared to the one in the HR simulation. Latitudinally,

394 AR distributions show hemispheric contrast, as compared to the LR, ARs in HR are
395 located more southward in the Pacific sector but more northward in the Atlantic sector.

396
397 Figure 7 shows AR characteristic distribution in the North Pacific for E3SM HR, LR and
398 ERA5. Apparently, E3SM produces more AR events than the reanalysis in nearly all
399 fields and across all scales. We also evaluated the precipitation associated with
400 landfalling ARs in California in both HR and LR simulations, as in Fig. 8. It is notable
401 that both models simulate systematically higher precipitation than ERA5 for all rainfall
402 intensity categories. It is also clear that the precipitation bias in HR simulation is larger
403 than LR simulation, except in the light rainfall ($< \sim 6\text{mm/day}$) category. Similarly, better
404 topographic representation in high resolution version of the model does not improve
405 precipitation simulation is also reported in Harrop et al. (2023), especially when the bias
406 in the low resolution model is substantially high.

407 408 **3.4. Sensitivity of AR characteristics to ARDT**

409
410 ARDTs are generally threshold-based, mostly using the intensity of moisture transport
411 with some geographical constraints that limit the AR spatial extent and some
412 geometrical constraints that preserve their nature as “long and narrow” filaments of
413 moisture. The different choices made by ARDT developers essentially amount to
414 different definitions of ARs (O’Brien et al. 2022), all of which are qualitatively consistent
415 with the definition in the AMS glossary (Ralph et al., 2018). For example, the Mundhenk
416 algorithm (Mundhenk et al. 2016) calculates integrated water vapor transport (IVT)
417 anomalies relative to the historical period and uses a fixed relative threshold to identify
418 ARs that are above a certain percentile of the historical simulation. The
419 TempestExtremes (TE; Ullrich et al. 2021) method, as another example, uses relative
420 threshold on the Laplacian of the IVT field rather than the IVT field itself.

421
422 In this application of the metrics package, we examine how ARs in ERA5 are sensitive
423 to the choice of ARDT. In addition to TE-based AR tags, we use AR tags detected using
424 the Mundhenk_v3 algorithm for comparison. Despite significant differences in their
425 associated algorithms, results from ARTMIP showed their performance was similar and
426 close to the mean among all ARDTs (Shields et al., 2018). Table 2 shows agreement of
427 landfalling ARs detected using these two ARDTs, as % values of IoU (AR concurrence
428 normalized by total occurrence of the ARs in both methods). The level of consistency
429 ranges from 56% to 83%, which suggests that TE and Mundhenk detect ARs
430 concurrently most of the time, but with asynchronous discrepancies, possibly at the
431 timing of the landfall and the end of the AR life cycle.

432

433 For AR characteristics over the oceans, the Mundhenk method detects larger ARs in
434 area, width, and length compared to TE (Fig. 9). Such differences are attributable to the
435 different thresholds for tagging the moisture field in the two ARDTs. The results
436 presented here are obtained from the default criteria, i.e. in TE, ARs are tagged when
437 the Laplacian of the IVT ≤ -20000 , while Mundhenk uses a static $250 \text{ kg m}^{-1} \text{ s}^{-1}$
438 threshold on the IVT field. We might expect different results by altering these threshold
439 numbers. ARs are also present at more northward latitudes with Mundhenk than TE as
440 indicated by positive biases. Zonally, AR distributions exhibit more hemispherical
441 contrast, with Mundhenk showing more westward located (positive biases) ARs in the
442 Pacific sector but more eastward located (negative biases) ARs in the Atlantic sector.
443 Apart from TE and Mundhenk, examples of AR geometry patterns from a few other
444 ARDTs are shown in supplementary Fig. S3, all showing the results from different
445 criteria for moisture tagging.

446

447 **3.5 Landfalling AR precipitation in CMIP5/6 models**

448

449 Apart from comparing AR properties, one useful capability of the ARMP is for analyzing
450 and quantifying any climate fields that are associated with ARs, e.g., precipitation, which
451 is an important indicator of the intensity of a landfalling AR. Here we evaluate landfalling
452 AR precipitation in the CMIP5 and CMIP6 models, with the ERA5 reanalysis and
453 MSWEP (Beck et al. 2017) gridded product as reference. Fig. 10 shows that compared
454 to the observations, landfalling precipitation differences in the models are generally
455 much larger than in reanalysis. The models show dry biases in most regions,
456 particularly large (up to -7.7 mm/day) in California, Pacific Northwest, Iceland and
457 Greenland.

458

459 As it is unclear if these biases are mainly due to general precipitation biases, or AR
460 activity bias, we further examine model precipitation bias diagnostics regardless of AR
461 activity (Fig. 11a) and AR frequency bias metrics (Fig. 11b) separately. For total
462 precipitation in the models, structural biases as in Fig. 10 are absent, but AR landfalls
463 are less frequent in the Pacific Northwest, Iceland, and Greenland. This suggests that
464 the systematic dry AR precipitation biases over these regions are primarily due to the
465 insufficient number of landfalling ARs in the models. For California, similar results do not
466 hold for all the models, for example, total precipitation in CCSM4 is 3.4 mm/day higher
467 than the reanalysis and AR landfalls are 6% more frequent, but the AR-related rainfall
468 has a dry bias of -0.5 mm/day . This suggests that landfalling ARs in CCSM4 are less
469 intense, suggesting a potential direction for model improvement.

470

471 **3.6 Landfalling AR peak day**

472

473 3.6.1 Comparison among reanalyses

474

475 Seasonality of AR landfalls is one of the important metrics for understanding AR
476 variability and impacts. Here we analyze landfalling AR seasonality over various regions
477 of the globe among three reanalysis products. We perform a Fourier transform on the
478 10-year long-term daily mean AR histogram to find its peak date based on the phase of
479 the first Fourier mode. Results indicate that the AR peak days agree well among
480 reanalyses for most regions, with small differences of only a few days. Large
481 discrepancies are noted for Australia and western Africa: In Australia, AR landfall peaks
482 nearly a month behind in JRA-55C than MERRA-2, while in west Africa, AR landfall in
483 MERRA-2 peaks 46 days behind ERA5.

484

485 Details of these differences are depicted in the histogram plots. For West Africa, AR
486 landfalls have two peaks in ERA5 and MERRA-2, one being in September, followed by
487 another peak in November. In ERA5, the peak in November is the main peak, while in
488 MERRA-2, the September peak is comparable to the November peak, resulting in an
489 earlier peak day from the Fourier phase spectrum. JRA-55C, in contrast, has only one
490 peak in November, and the AR landfall event counts are fewer than the other two
491 products over the entire year, indicative of smaller year to year variability.

492

493 Seasonal distribution of AR landfalls in Australia in the three reanalyses exhibit similar
494 differences to those in western Africa. In ERA5 and MERRA-2, there are two peaks in
495 February and June, but only one peak presents in JRA-55C in June. This explains the
496 relative late peak day in JRA-55C. While the main peak in ERA5 is in June, in MERRA-
497 2, the main peak is in February, which is consistent with the metrics result that MERRA-
498 2 has the earliest peak day. Similarly, the JRA-55C has a smaller number of landfalling
499 ARs, although the interannual variability is comparable to the other two reanalyses.

500

501 3.6.2 AR seasonality in climate models

502

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

509

510 In regions like South America, Baja, UK and Western Europe, the models show
511 systematic late peak biases, and in South Africa, AR peaks earlier than the reanalyses.
512 The exact cause of these structural biases in the models is likely indicative of persistent

513 and ubiquitous timing issues in the shift of the storm track that is common among
514 models. It is worth noting that the model biases in the West Africa region are
515 significantly larger than other regions, with peak day difference up to 6 months as
516 compared to the reanalysis. Looking at the AR counts histograms over the course of the
517 year in this region in the CCSM3 and MRI-ESM2-0 models (Fig. 14), it is clear that AR
518 landfall seasonality in both models is completely out of phase with ERA5. This is
519 especially true for the MRI-ESM2-0 model, where AR landfall peaks in June, which is in
520 opposition to the climatology in ERA5. The large discrepancy is probably because of the
521 large spread in the atmospheric circulations in climate models over the West Africa
522 region, as large spread among CMIP5/6 models in capturing atmospheric dynamic
523 responses (Monerie et al. 2020), the lack of jet-rainfall coupling (Whittleston et al. 2017),
524 and bias in simulating mesoscale convective systems (Jenkins et al. 2002) in climate
525 models are noted. Although high resolution regional modeling may be capable of
526 improving rainfall in this region (Sylla et al. 2009), the dynamics-rainfall coupling does
527 not appear to be improved in high resolution global models such as the E3SM (Caldwell
528 et al. 2019; Golaz et al. 2019). Therefore, challenges remain in modeling the AR water
529 cycle in west Africa.

530

531 **4. Summary and discussion**

532

533 In this study we have introduced a metrics framework, namely ARMP, for the objective
534 evaluation of ARs in climate models and reanalysis, and illustrated the potential for its
535 use with five example case-studies to illustrate the scope of potential applications. The
536 metrics-based analyses are designed for systematic diagnosis of AR biases in climate
537 models. In our example application applying the package to CMIP5 and CMIP6 models,
538 we have shown that AR tracks in the south Atlantic are positioned farther poleward
539 compared to the ERA5 reanalysis, while in the south Pacific, tracks are biased towards
540 the equator. Over western Africa, we found that most climate models do a poor job at
541 capturing the AR peak season, while it is generally consistent among reanalyses.

542

543 In the application of comparing AR characteristics represented in high- and low-
544 resolution model simulation, while biases are not generally reduced in high-resolution
545 configuration, substantial differences are noted between the two simulations. For
546 example, in the North Atlantic basin, AR tracks in the E3SM-HR are shifted more
547 northward than in the E3SM-LR simulation. In addition to model evaluation and model
548 and reanalysis intercomparison, we have shown how our metrics package can be used
549 to identify structural differences resulting from the choice of AR detector (ARDT). For
550 instance, we demonstrated that ARs detected with the Mundhenk method are
551 systematically larger in size, width and length compared to TE.

552

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

561

562 Our metrics package assembles a suite of established and newly introduced AR metrics
563 into one framework, facilitating objective evaluation of ARs with a diverse suite of input
564 data, as well as intercomparison of ARs as simulated by multiple climate models. These
565 metrics can be routinely applied for model benchmarking and during development
566 cycles to monitor changes in AR characteristics across model versions or generations
567 and be used to set objective targets for the improvement of models. One expected
568 application is the routine benchmarking of AR in simulations with increasingly higher
569 resolution models. More frequent metrics evaluation of simulated ARs such as this
570 could further our understanding of model bias and error characteristics, and potentially
571 assist developers in making choices associated with new model versions. Furthermore,
572 it can provide a quantitative measure for operational centres to perform near real-time
573 climate and extreme events impact assessment along with their forecast cycles, which
574 can facilitate their decision-making process.

575

576 The collection of metrics included in our metrics package will be augmented to gauge
577 additional AR characteristics. At the time of the submission of this manuscript, it is being
578 configured to be a part of the PMP. Looking forward, we welcome community
579 contributions to successive development of the package. Inspired by Zarzycki et al.
580 (2021), there is also a potential that these metrics can be applied for research beyond
581 ARs, such as mesoscale meteorological features, regional hydrological extremes such
582 as floods and droughts, and large-scale climate modes.

583

584

585 **Appendix**

586

587 This section includes mathematical expressions of commonly used model evaluation
588 metrics.

589

590 **A1. Mean bias**

591

592 The mean bias is mathematically expressed as

593

594

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

595

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

597

598

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

599

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

601

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

604

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

605

606

607 where \bar{x}_i is sample arithmetic mean, μ_i is population mean, s_i is sample variance, and
608 n_i is sample size. A positive z-score indicates that the value is above the mean. The
609 higher the z-score, the further above the mean the value is, and vice versa. A result is
610 considered statistically significant at the 95% confidence level if the magnitude of the z-
611 score is greater than 1.96.

612

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

617

618 A2. Spatial pattern similarity

619

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

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

623

624 where, x_i and y_i are the values of the two spatial patterns at location i (or grid point i in
625 gridded data product), \bar{x} and \bar{y} are the means of the values of the two patterns, and n is

626 the total number of locations. This equation essentially measures the degree to which
627 the values of the two spatial patterns vary together. If they vary together perfectly, r will
628 be 1. If they vary together inversely, r will be -1. If there's no linear relationship between
629 the patterns, r will be 0.

630
631 The statistical significance of correlation is determined by the two-tailed p-value of the
632 cumulative distribution function (CDF) of the t-statistic, as

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

635
636 The the t-statistic t is given by

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

637
638 where r is the correlation coefficient, and n_e is the effective sample size. Although there
639 are a number of methods to estimate the effective geographic sample size (e.g., Griffith
640 2013).

641
642
643

644 **Acknowledgment**

645
646 This work is performed under the auspices of the U.S. Department of Energy (DOE) by
647 Lawrence Livermore National Laboratory (LLNL) under Contract No. DE-AC52-
648 07NA27344 and is mainly supported by the Regional and Global Model Analysis
649 (RGMA) program of the U.S. DOE Office of Science (OS) Biological and Environmental
650 Research (BER) program. This material is based upon work supported by the U.S.
651 Department of Energy, Office of Science, Office of Biological and Environmental
652 Research, Climate and Environmental Sciences Division, Regional & Global Model
653 Analysis Program. Resources of the National Energy Research Scientific Computing
654 Center (NERSC) were used. The research was partially supported by the Office of
655 Science of the U.S. Department of Energy under Contract Number DE-AC02-
656 05CH11231 and under award Number DE-SC0023519. This research was also
657 supported in part by the Environmental Resilience Institute, funded by Indiana
658 University's Prepared for Environmental Change Grand Challenge initiative and in part
659 by Lilly Endowment, Inc., through its support for the Indiana University Pervasive
660 Technology Institute. We acknowledge the World Climate Research Programme, which,
661 through its Working Group on Coupled Modeling, coordinated and promoted CMIP6.
662 We thank the climate modeling groups for producing and making available their model
663 output, the Earth System Grid Federation (ESGF) for archiving the data and providing

664 access, and the multiple funding agencies that support CMIP6 and ESGF. The authors
665 acknowledge Antony Hoang and Ana Ordonez for computing and technical support, and
666 Christine Shields, Yang Zhou and Allison Collow for their help on data and discussion.
667

668 **Code and data availability**

669 The ARMP code is hosted on GitHub <https://github.com/PCMDI/ARMP>. The initial
670 release is also available on Zenodo DOI 10.5281/zenodo.14188789. Users are strongly
671 recommended to download the source code from GitHub to ensure access to the latest
672 changes, updates and improvements of the package.
673

674 **Author contribution**

675 BD implemented the codes and developed the diagnostic results. All authors
676 contributed to the writing of the manuscript.
677

678 **Competing interests**

679 At least one of the (co-)authors is a member of the editorial board of Geoscientific
680 Model Development.
681

682 **References**

- 683
- 684 Algarra, I., Nieto, R., Ramos, A. M., Eiras-Barca, J., Trigo, R. M., & Gimeno, L. (2020).
685 Significant increase of global anomalous moisture uptake feeding land-falling
686 atmospheric rivers. *Nature communications*, 11 (1), 5082.
- 687 Beck, H. E., Van Dijk, A. I., Levizzani, V., Schellekens, J., Miralles, D. G., Martens, B., &
688 De Roo, A. (2017). MSWEP: 3-hourly 0.25 global gridded precipitation (1979–2015)
689 by merging gauge, satellite, and reanalysis data. *Hydrology and Earth System
690 Sciences*, 21(1), 589-615.
- 691 Buizza, R., Poli, P., Rixen, M., Alonso-Balmaseda, M., Bosilovich, M. G., Brönnimann,
692 S., ... & Vasselali, A. (2018). Advancing global and regional reanalyses. *Bulletin of
693 the American Meteorological Society*, 99(8), ES139-ES144.
- 694 Caldwell, P. M., Mamejtanov, A., Tang, Q., Van Roekel, L. P., Golaz, J. C., Lin, W., ... &
695 Zhou, T. (2019). The DOE E3SM coupled model version 1: Description and results
696 at high resolution. *Journal of Advances in Modeling Earth Systems*, 11(12), 4095-
697 4146.
- 698 Chapman, W. E., Subramanian, A. C., Delle Monache, L., Xie, S. P., & Ralph, F. M.
699 (2019). Improving atmospheric river forecasts with machine learning. *Geophysical
700 Research Letters*, 46(17-18), 10627-10635.
- 701 DeFlorio, M. J., Waliser, D. E., Guan, B., Lavers, D. A., Ralph, F. M., & Vitart, F. (2018).
702 Global assessment of atmospheric river prediction skill. *Journal of
703 Hydrometeorology*, 19(2), 409-426.

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

706 Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., & Taylor, K.
707 E. (2016). Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6)
708 experimental design and organization. *Geoscientific Model Development*, 9(5),
709 1937-1958.

710 Gelaro, R., McCarty, W., Suárez, M. J., Todling, R., Molod, A., Takacs, L., . . . others
711 (2017). The modern-era retrospective analysis for research and applications,
712 version 2 (MERRA-2). *Journal of climate*, 30 (14), 5419–5454.

713 Gershunov, A., Shulgina, T., Clemesha, R. E., Guirguis, K., Pierce, D. W., Dettinger, M.
714 D., . . . others (2019). Precipitation regime change in western North America: the role
715 of atmospheric rivers. *Scientific reports*, 9 (1), 9944.

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

718 Gleckler, P. J., Taylor, K. E., & Doutriaux, C. (2008). Performance metrics for climate
719 models. *Journal of Geophysical Research: Atmospheres*, 113(D6).

720 Golaz, J. C., Caldwell, P. M., Van Roekel, L. P., Petersen, M. R., Tang, Q., Wolfe, J. D., ...
721 & Zhu, Q. (2019). The DOE E3SM coupled model version 1: Overview and
722 evaluation at standard resolution. *Journal of Advances in Modeling Earth Systems*,
723 11(7), 2089-2129.

724 Griffith, D. A. (2013). Establishing qualitative geographic sample size in the presence of
725 spatial autocorrelation. *Annals of the Association of American Geographers*, 103(5),
726 1107-1122.

727 Guan, B., Molotch, N. P., Waliser, D. E., Fetzer, E. J., & Neiman, P. J. (2010). Extreme
728 snowfall events linked to atmospheric rivers and surface air temperature via satellite
729 measurements. *Geophysical Research Letters*, 37 (20).

730 Guan, B., & Waliser, D. E. (2017). Atmospheric rivers in 20 year weather and climate
731 simulations: A multimodel, global evaluation. *Journal of Geophysical Research:*
732 *Atmospheres*, 122(11), 5556-5581.

733 Guan, B., Waliser, D. E., & Ralph, F. M. (2023). Global application of the atmospheric
734 river scale. *Journal of Geophysical Research: Atmospheres*, 128(3),
735 e2022JD037180.

736 Harrop, B., Leung, L. and Ullrich P. (2023). Improving Simulations of Atmospheric Rivers
737 and Heat Waves in the Coupled E3SM. FY2023 First Quarter Performance Metric.
738 DOE/SC-CM-23-001

739 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., . . .
740 others (2020). The ERA5 global reanalysis. *Quarterly Journal of the Royal*
741 *Meteorological Society*, 146 (730), 1999–2049.

742 Hoyer, S. & Hamman, J., (2017). xarray: N-D labeled Arrays and Datasets in Python.
743 Journal of Open Research Software. 5(1), p.10. DOI:
744 <https://doi.org/10.5334/jors.148>

745 Huang, X., Swain, D. L., & Hall, A. D. (2020). Future precipitation increase from very high
746 resolution ensemble downscaling of extreme atmospheric river storms in California.
747 Science advances, 6(29), eaba1323.

748 Huang, J., Zhang, C., & Prospero, J. M. (2009). African aerosol and large-scale
749 precipitation variability over West Africa. *Environmental Research Letters*, 4(1),
750 015006.

751 Hui, W. J., Cook, B. I., Ravi, S., Fuentes, J. D., & D'Odorico, P. (2008). Dust-rainfall
752 feedbacks in the West African Sahel. *Water Resources Research*, 44(5).

753 Inda-Díaz, H. A., O'Brien, T. A., Zhou, Y., & Collins, W. D. (2021). Constraining and
754 characterizing the size of atmospheric rivers: A perspective independent from the
755 detection algorithm. *Journal of Geophysical Research: Atmospheres*, 126(16),
756 e2020JD033746.

757 Kobayashi, S., Ota, Y., Harada, Y., Ebita, A., Moriya, M., Onoda, H., . . . others (2015).
758 The JRA-55 reanalysis: General specifications and basic characteristics. *Journal of*
759 *the Meteorological Society of Japan*. Ser. II , 93 (1), 5–48.

760 Lee, J., P. J. Gleckler, M.-S. Ahn, A. Ordonez, P. Ullrich, K. R. Sperber, K. E. Taylor, Y.
761 Y. Planton, E. Guilyardi, P. Durack, C. Bonfils, M. D. Zelinka, L.-W. Chao, B. Dong,
762 C. Doutriaux, C. Zhang, T. Vo, J. Boutte, M. F. Wehner, A. G. Pendergrass, D. Kim,
763 Z. Xue, A. T. Wittenberg, and J. Krasting, (2024): Systematic and Objective
764 Evaluation of Earth System Models: PCMDI Metrics Package (PMP) version 3.
765 Geoscientific Model Development, 17, 3919–3948, doi: 10.5194/gmd-17-3919-
766 2024.

767

768 Leung, L. R., & Qian, Y. (2009). Atmospheric rivers induced heavy precipitation and
769 flooding in the western us simulated by the WRF regional climate model.
770 Geophysical research letters, 36 (3).

771 Lora, J. M., Mitchell, J. L., Risi, C., & Tripathi, A. E. (2017). North pacific atmospheric rivers
772 and their influence on western north America at the last glacial maximum.
773 Geophysical Research Letters, 44 (2), 1051–1059.

774 Massoud, E. C., Espinoza, V., Guan, B., & Waliser, D. E. (2019). Global Climate Model
775 Ensemble Approaches for Future Projections of Atmospheric Rivers. *Earth's Future*,
776 7: 1136–11511151.

777 Mundhenk, B. D., Barnes, E. A., & Maloney, E. D. (2016). All-season climatology and
778 variability of atmospheric river frequencies over the north pacific. *Journal of Climate*,
779 29 (13), 4885–4903.

780 Nardi, K. M., Barnes, E. A., & Ralph, F. M. (2018). Assessment of numerical weather
781 prediction model reforecasts of the occurrence, intensity, and location of

782 atmospheric rivers along the West Coast of North America. *Monthly Weather*
783 *Review*, 146(10), 3343-3362.

784 Neiman, P. J., Ralph, F. M., Moore, B. J., Hughes, M., Mahoney, K. M., Cordeira, J. M.,
785 & Dettinger, M. D. (2013). The landfall and inland penetration of a flood-producing
786 atmospheric river in arizona. part i: Observed synoptic-scale, orographic, and
787 hydrometeorological characteristics. *Journal of Hydrometeorology*, 14 (2), 460–484.

788 Neiman, P. J., Schick, L. J., Ralph, F. M., Hughes, M., & Wick, G. A. (2011). Flooding in
789 western washington: The connection to atmospheric rivers. *Journal of*
790 *Hydrometeorology*, 12 (6), 1337–1358.

791 Neiman, P. J., White, A. B., Ralph, F. M., Gottas, D. J., & Gutman, S. I. (2009). A water
792 vapour flux tool for precipitation forecasting. In *Proceedings of the institution of civil*
793 *engineers-water management* (Vol. 162, pp. 83–94).

794 Newman, M., Kiladis, G. N., Weickmann, K. M., Ralph, F. M., & Sardeshmukh, P. D.
795 (2012). Relative contributions of synoptic and low-frequency eddies to time-mean
796 atmospheric moisture transport, including the role of atmospheric rivers. *Journal of*
797 *climate*, 25(21), 7341-7361.

798 O'Brien, T. A., Risser, M. D., Loring, B., Elbashandy, A. A., Krishnan, H., Johnson, J., &
799 Collins, W. D. (2020). Detection of atmospheric rivers with inline uncertainty
800 quantification: TECA-BARD v1.0.1. *Geoscientific Model Development*, 13(12),
801 6131–6148. <https://doi.org/10.5194/gmd-13-6131-2020>

802 O'Brien, T., Wehner, M., Payne, A., Shields, C., Rutz, J., Leung, L.-R., . . . others (2022).
803 Increases in Future AR Count and Size: Overview of the ARTMIP Tier 2 CMIP5/6
804 Experiment. *Journal of Geophysical Research: Atmospheres*, 127 (6).

805 Payne, A. E., Demory, M. E., Leung, L. R., Ramos, A. M., Shields, C. A., Rutz, J. J., ... &
806 Ralph, F. M. (2020). Responses and impacts of atmospheric rivers to climate
807 change. *Nature Reviews Earth & Environment*, 1(3), 143-157.

808 Ramachandran, J., & Aschheim, M. A. (2005). Sample size and error in the determination
809 of mode shapes by principal components analysis. *Engineering structures*, 27(14),
810 1951-1967.

811 Ralph, F., Coleman, T., Neiman, P., Zamora, R., & Dettinger, M. (2013). Observed
812 impacts of duration and seasonality of atmospheric-river landfalls on soil moisture
813 and runoff in coastal northern California. *Journal of Hydrometeorology*, 14 (2), 443–
814 459.

815 Ralph, F.M. and Dettinger, M. (2011). Storms, floods, and the science of atmospheric
816 rivers. *Eos, Transactions American Geophysical Union*, 92(32), 265-266.

817 Ralph, F. M., Neiman, P. J., Kiladis, G. N., Weickmann, K., & Reynolds, D. W. (2011). A
818 multiscale observational case study of a pacific atmospheric river exhibiting
819 tropical—extratropical connections and a mesoscale frontal wave. *Monthly Weather*
820 *Review* , 139 (4), 1169–1189.

821 Ralph, F. M., Neiman, P. J., Wick, G. A., Gutman, S. I., Dettinger, M. D., Cayan, D. R., &
822 White, A. B. (2006). Flooding on California's Russian River: Role of atmospheric
823 rivers. *Geophysical Research Letters*, 33 (13).

824 Ralph, F.M., Dettinger, M.D., Cairns, M.M., Galarneau, T.J. and Eylander, J. (2018).
825 Defining "atmospheric river": How the Glossary of Meteorology helped resolve a
826 debate. *Bulletin of the American Meteorological Society*, 99(4), pp.837-839.

827 Rutz, J. J., Steenburgh, W. J., & Ralph, F. M. (2014). Climatological characteristics of
828 atmospheric rivers and their inland penetration over the western United states.
829 *Monthly Weather Review* , 142 (2), 905–921.

830 Rutz, J.J., Shields, C.A., Lora, J.M., Payne, A.E., Guan, B., Ullrich, P., O'brien, T., Leung,
831 L.R., Ralph, F.M., Wehner, M. and Brands, S. (2019). The atmospheric river tracking
832 method intercomparison project (ARTMIP): Quantifying uncertainties in atmospheric
833 river climatology. *Journal of Geophysical Research: Atmospheres*, 124(24),
834 pp.13777-13802.

835 Shearer, Eric J.; Nguyen, Phu; Sellars, Scott L.; Analui, Bitu; Kawzenuk, Brian; Hsu, Kuo-
836 Lin; Sorooshian, Soroosh (2020). The Atmospheric River-CONNected objECT (AR-
837 CONNECT) algorithm applied to the National Aeronautics and Space Administration
838 (NASA) Modern-Era Retrospective Analysis for Research and Applications, Version
839 2 (MERRA V2) - 1983 to 2016. UC San Diego Library Digital Collections.
840 <https://doi.org/10.6075/J0D21W00>

841 Shields, C. A., Rosenbloom, N., Bates, S., Hannay, C., Hu, A., Payne, A. E., . . . Truesdale,
842 J. (2019). Meridional heat transport during atmospheric rivers in high-resolution
843 cesm climate projections. *Geophysical Research Letters*, 46 (24), 14702–14712.

844 Shields, C. A., Rutz, J. J., Leung, L. R., Ralph, F. M., Wehner, M., O'Brien, T., & Pierce,
845 R. (2019). Defining uncertainties through comparison of atmospheric river tracking
846 methods. *Bulletin of the American Meteorological Society*, 100 (2), ES93–ES96.

847 Solmon, F., Mallet, M., Elguindi, N., Giorgi, F., Zakey, A., & Konaré, A. (2008). Dust
848 aerosol impact on regional precipitation over western Africa, mechanisms and
849 sensitivity to absorption properties. *Geophysical Research Letters*, 35(24).

850 Taylor, K. E. (2001). Summarizing multiple aspects of model performance in a single
851 diagram. *Journal of geophysical research: atmospheres*, 106(D7), 7183-7192.

852 Ullrich, P. A., & Zarzycki, C. M. (2017). Tempestextremes: A framework for scale-
853 insensitive pointwise feature tracking on unstructured grids. *Geoscientific Model*
854 *Development*, 10 (3), 1069–1090.

855 Ullrich, P. A., Zarzycki, C. M., McClenny, E. E., Pinheiro, M. C., Stansfield, A. M., & Reed,
856 K. A. (2021). Tempestextremes v2. 1: A community framework for feature detection,
857 tracking and analysis in large datasets. *Geoscientific model development*
858 *discussions*, 2021 , 1–37.

859 Vo et al., (2024). xCDAT: A Python Package for Simple and Robust Analysis of Climate
860 Data. *Journal of Open Source Software*, 9(98), 6426,
861 <https://doi.org/10.21105/joss.06426>
862 Wilks, D. S. (2011). *Statistical methods in the atmospheric sciences* (Vol. 100). Academic
863 press.
864 Zarzycki, C. M., Ullrich, P. A., & Reed, K. A. (2021). Metrics for evaluating tropical
865 cyclones in climate data. *Journal of Applied Meteorology and Climatology*, 60 (5),
866 643–660.
867 Zhang, L., Zhao, Y., Cheng, T. F., & Lu, M. (2024). Future changes in global atmospheric
868 rivers projected by CMIP6 models. *Journal of Geophysical Research: Atmospheres*,
869 129, e2023JD039359
870 Zhao, A., Ryder, C. L., & Wilcox, L. J. (2022). How well do the CMIP6 models
871 simulate dust aerosols?. *Atmospheric Chemistry and Physics*, 22(3), 2095-2119.
872 Zhu, Y., & Newell, R. E. (1998). A proposed algorithm for moisture fluxes from
873 atmospheric rivers. *Monthly weather review* , 126 (3), 725–735
874
875
876

877 **Supplementary information**

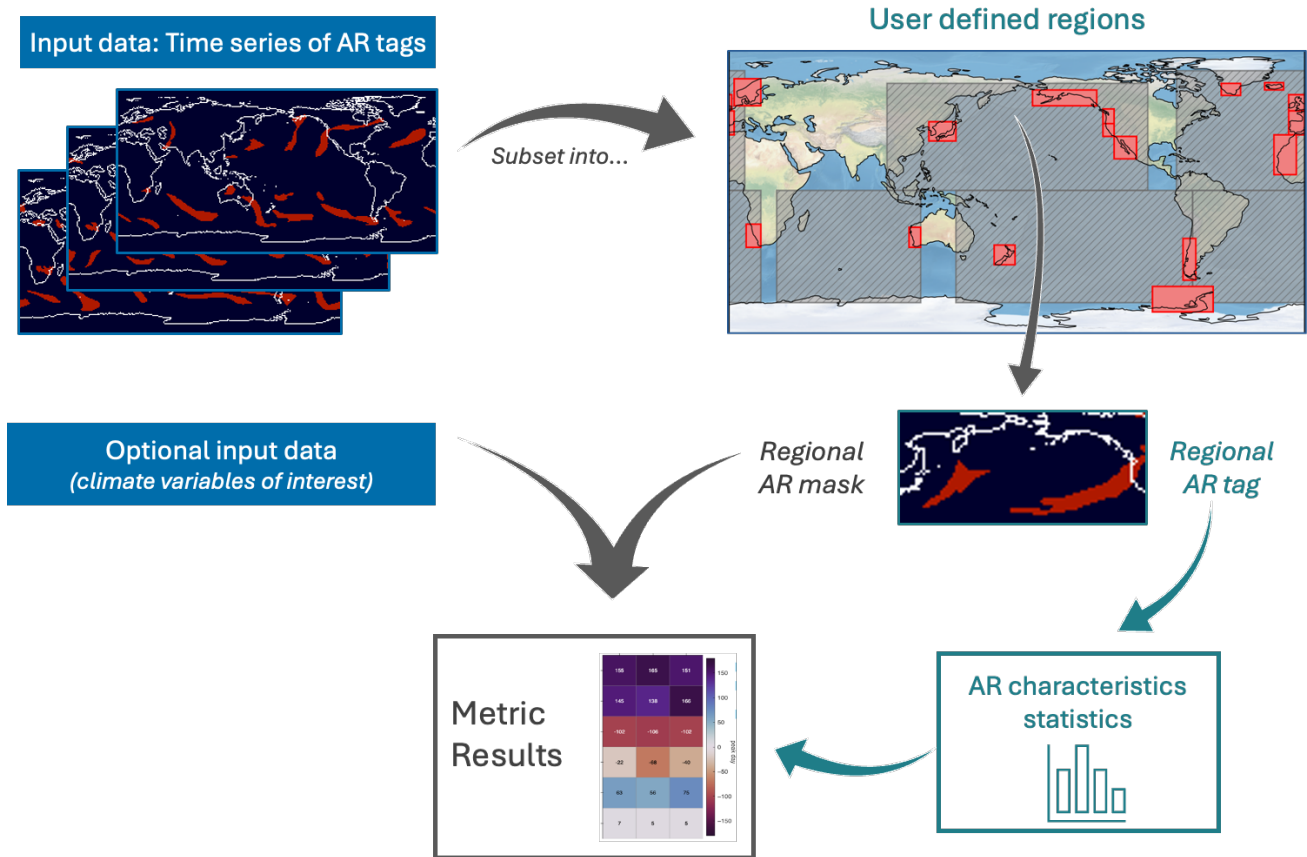
878 In a separate document

879

880

881 **Figures and tables**

882



883

884

885 Fig. 1. AR metric tool workflow. Input data include time slices of AR
 886 tags from ARDTs of user choice, and optional climate data
 887 associated with ARs. The data are then subset into user-defined
 888 rectangular domains (blue boxes for ocean basins, red boxes for
 889 landfall regions) for regional tags and masks. User preferred
 890 statistical tools are applied on the regional AR tags to obtain AR
 891 characteristics. Finally, AR characteristics and AR masked climate
 892 data are presented as metric results.

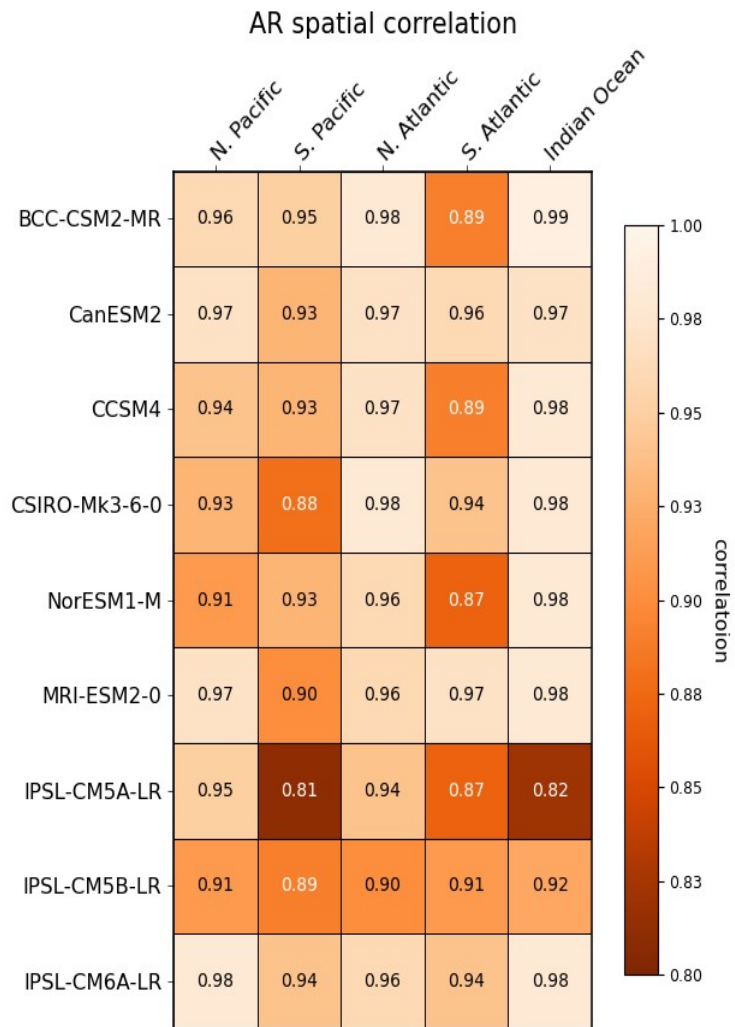
893

894
 895
 896
 897
 898
 899
 900

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.

	ARs over Ocean Basins						Landfalling ARs		
<u>metrics/diagnostics</u>	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		

901



902

903 Fig. 2. Spatial pattern correlation of AR frequency for the period 1979-
 904 1989 between ERA5 and climate models for major ocean basins.

905

906

907

908
909
910
911
912
913
914
915
916

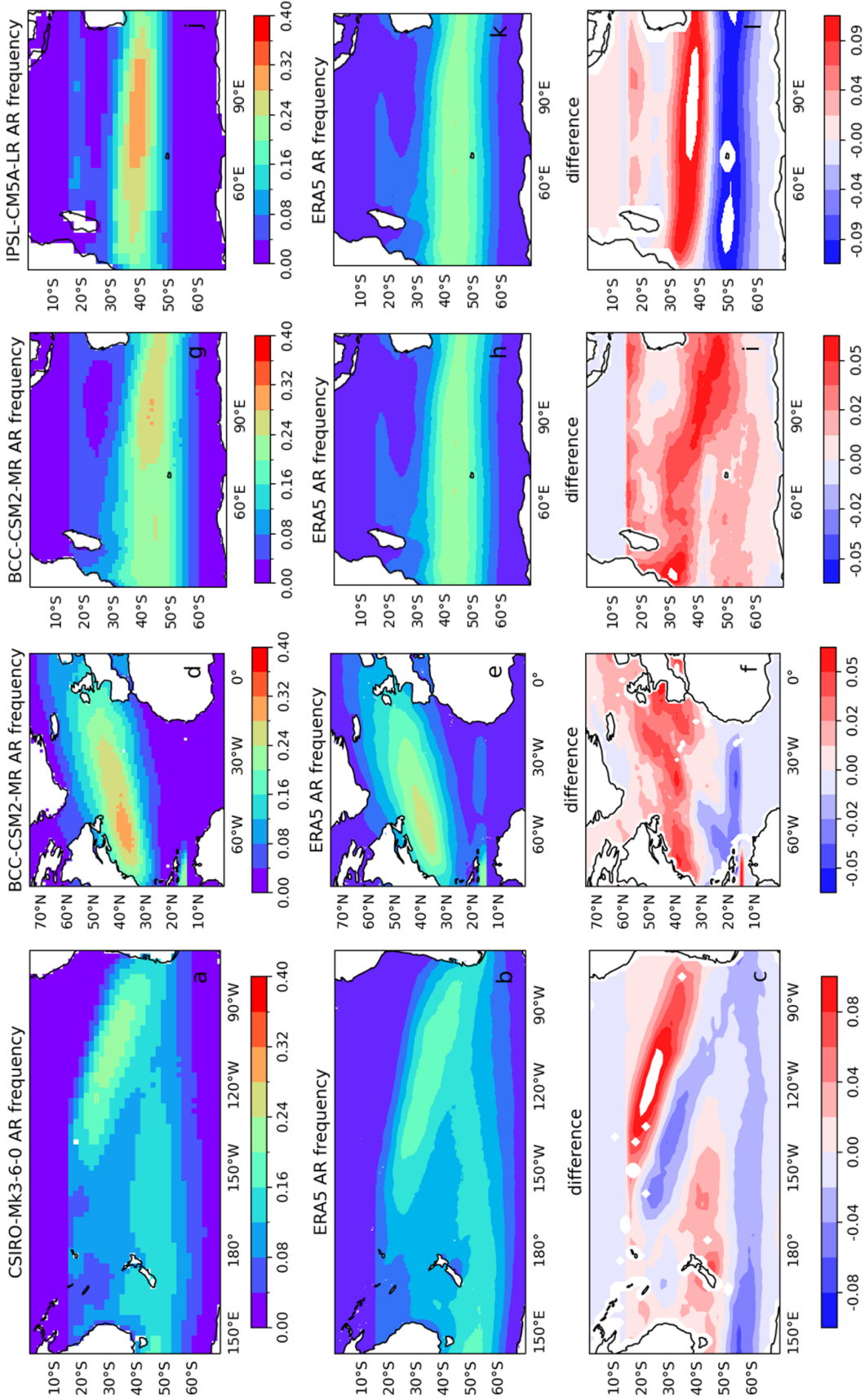


Fig. 3. AR frequency in the South Pacific for (a) CSIRO-Mk3-6-0, (b) ERA5 and their difference (c) as (a) - (b). (d)-(f), (g)-(i), (j)-(l) are same as (a)-(c) but for AR frequency in the North Atlantic for BCC-CSM2-MR, Indian Ocean for BCC-CSM2-MR, and Indian ocean for IPSL-CM5A-LR respectively.

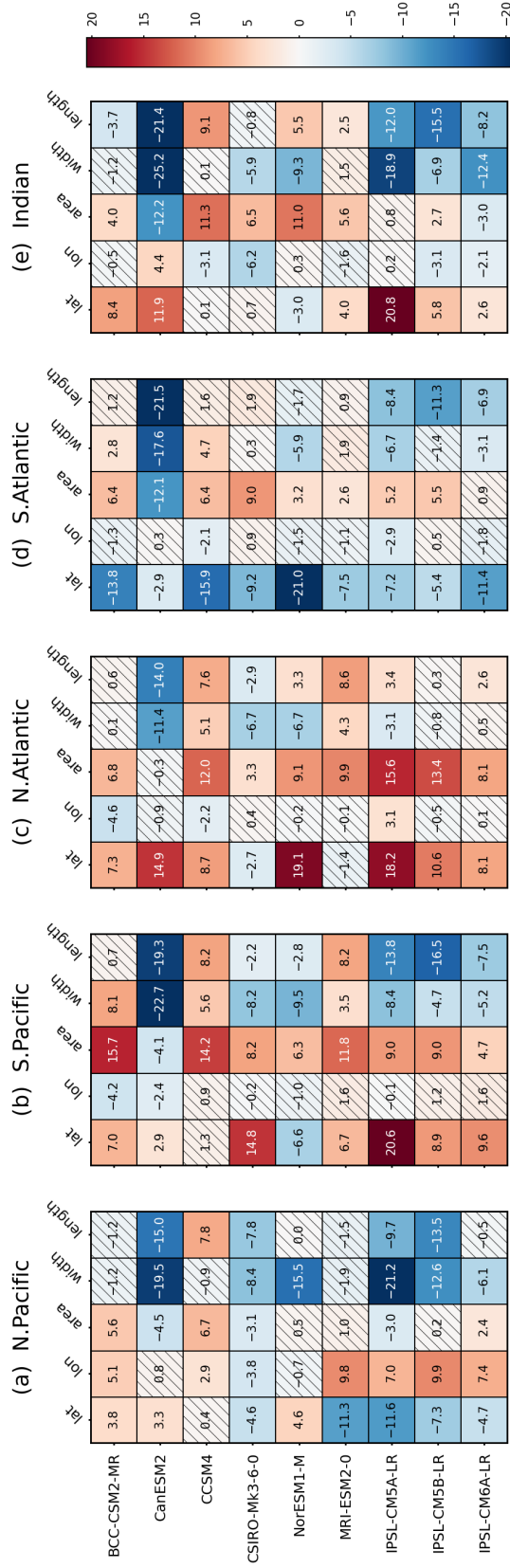
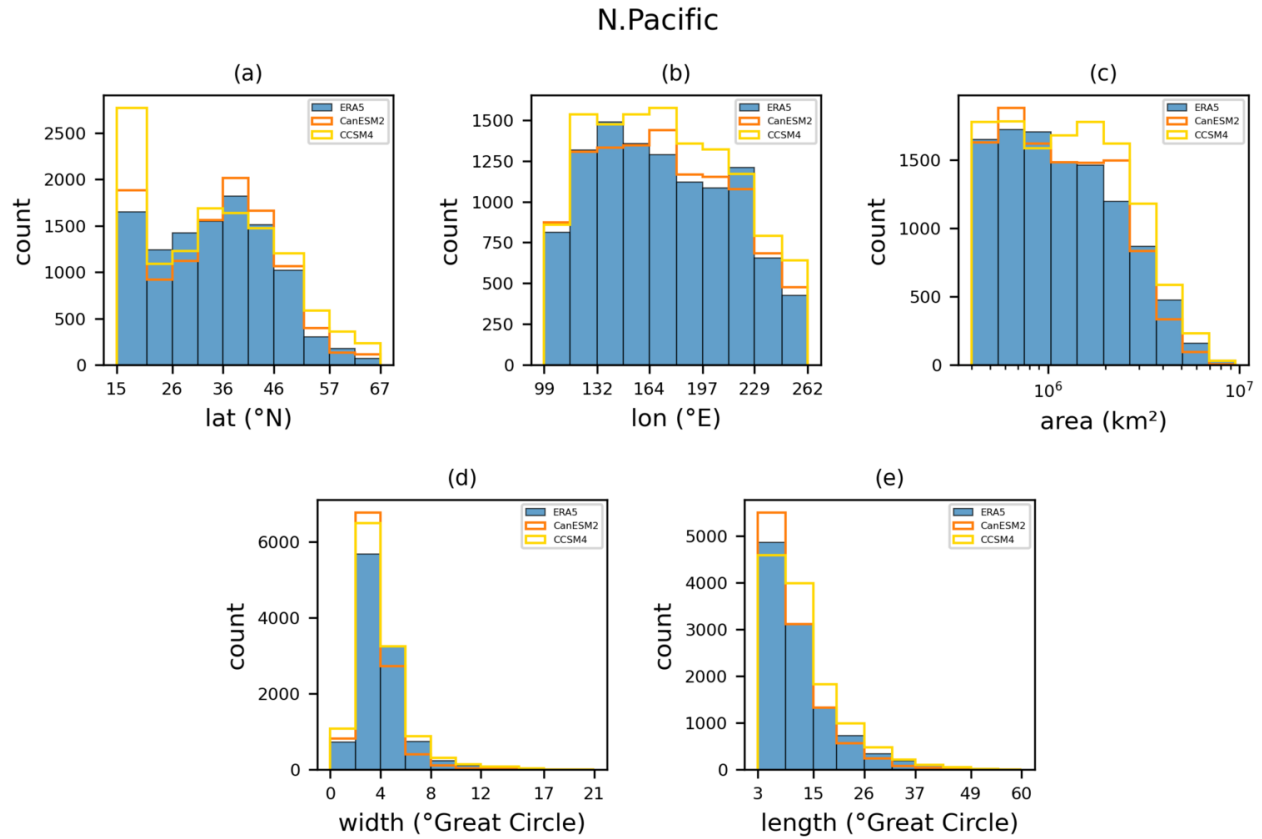


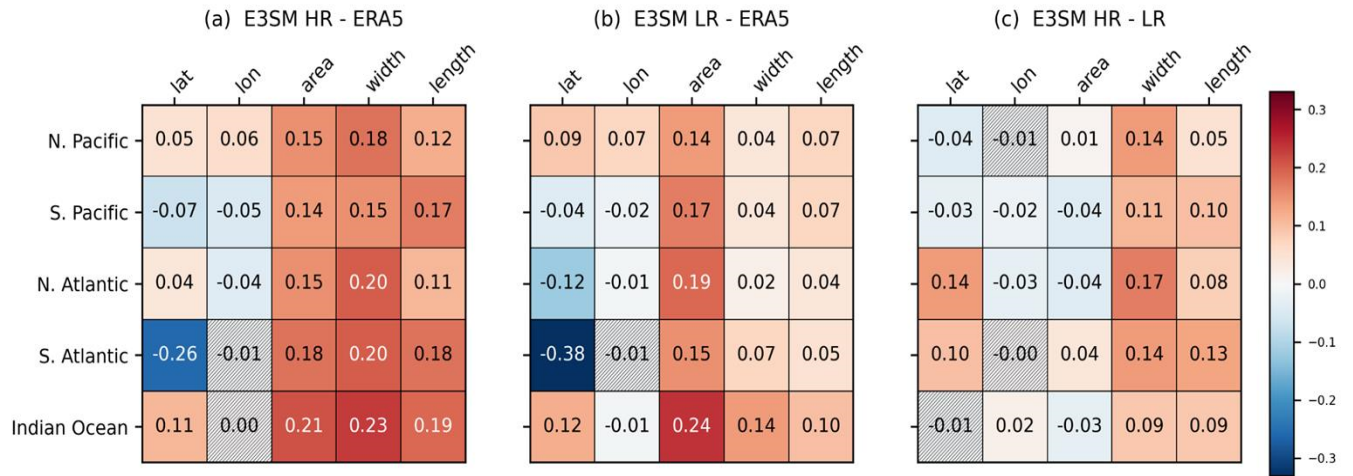
Fig. 4 AR characteristics bias (normalized as Z-score) in climate models for major ocean basins. Hatching indicates that the differences are statistically insignificant.



919

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

923
924



925
926
927
928
929
930
931
932

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.

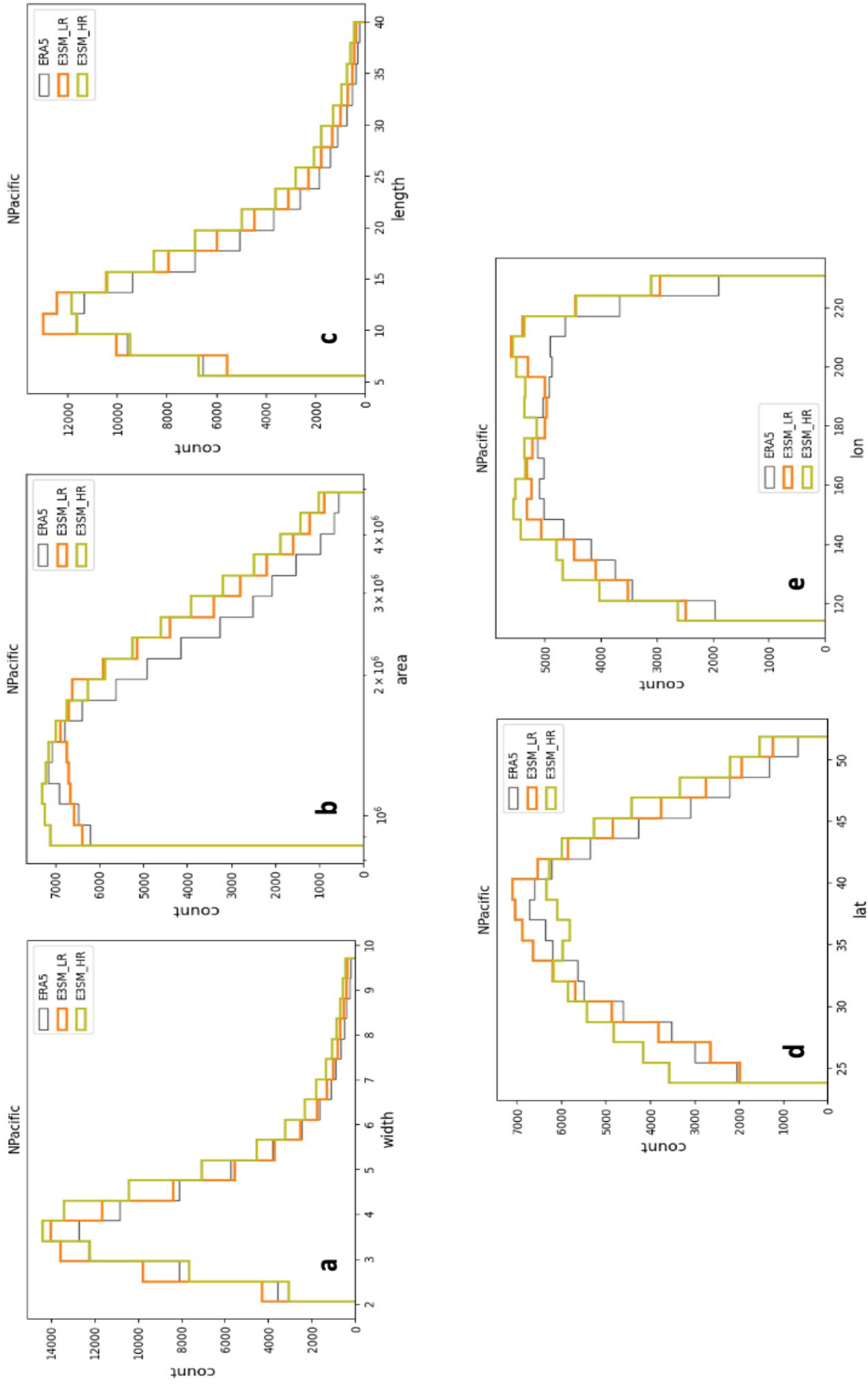
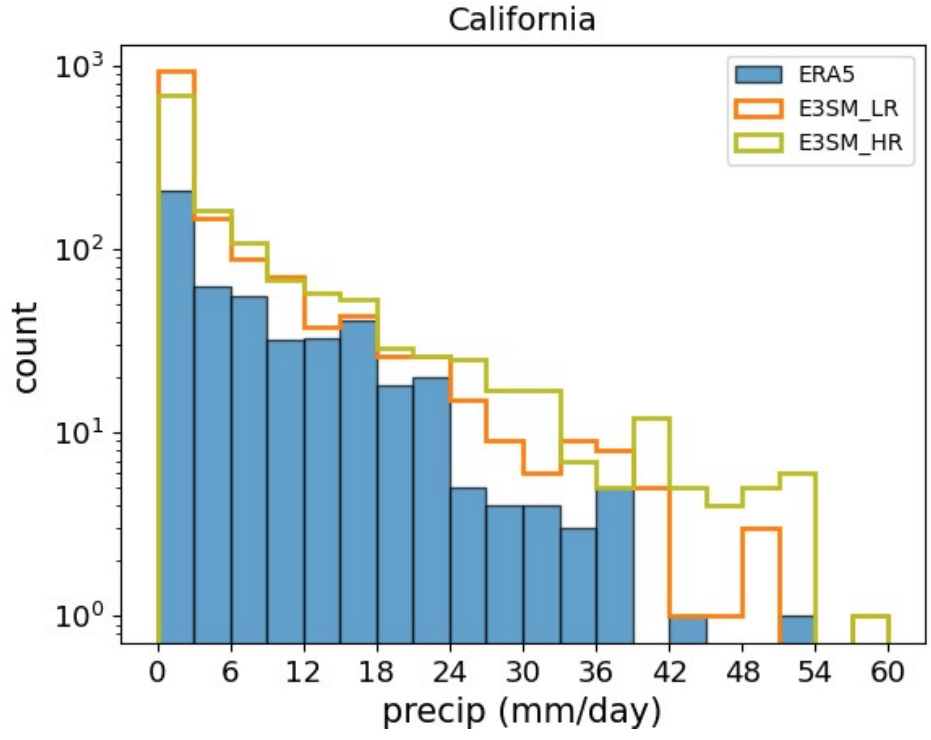


Fig. 7. AR characteristics distribution of (a) width (° great circle), (b) area (km²), (c) length (° great circle), (d) central latitude (°N), and (e) central longitude (°E) in the North Pacific for ERA5, E3SM LR and LR simulations.

934
935
936
937
938



939
940
941
942
943
944
945

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

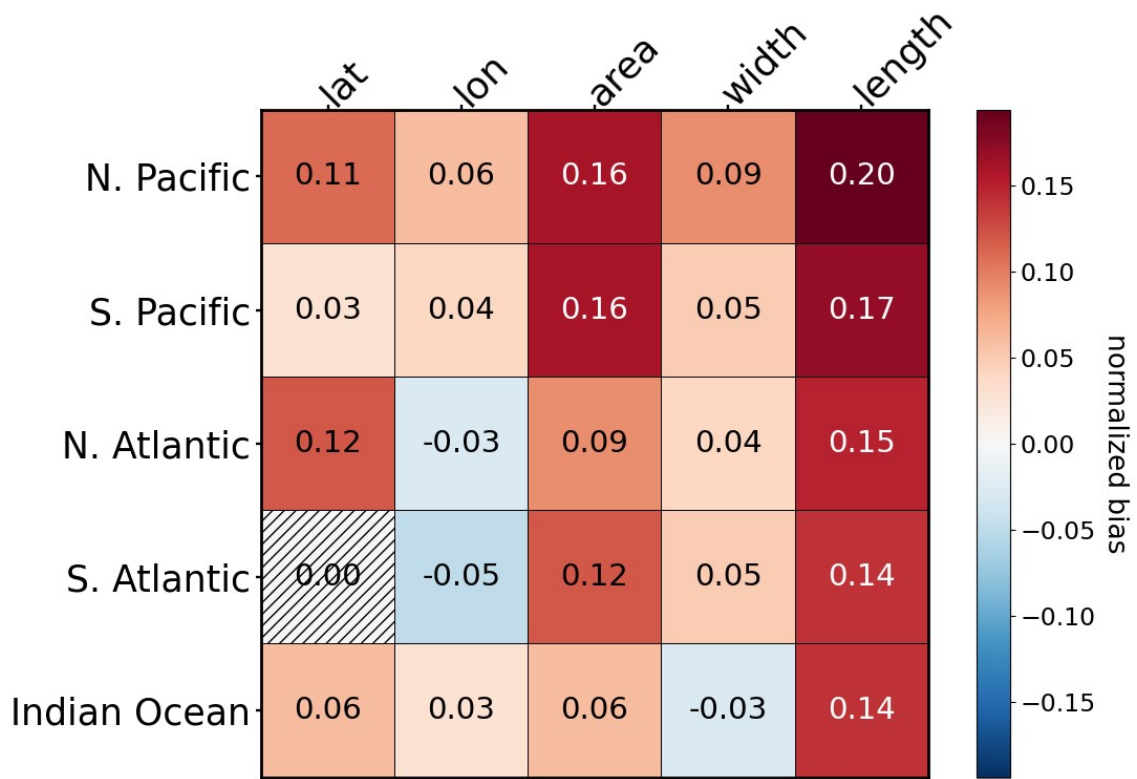
946
947
948
949
950
951
952
953
954

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

955

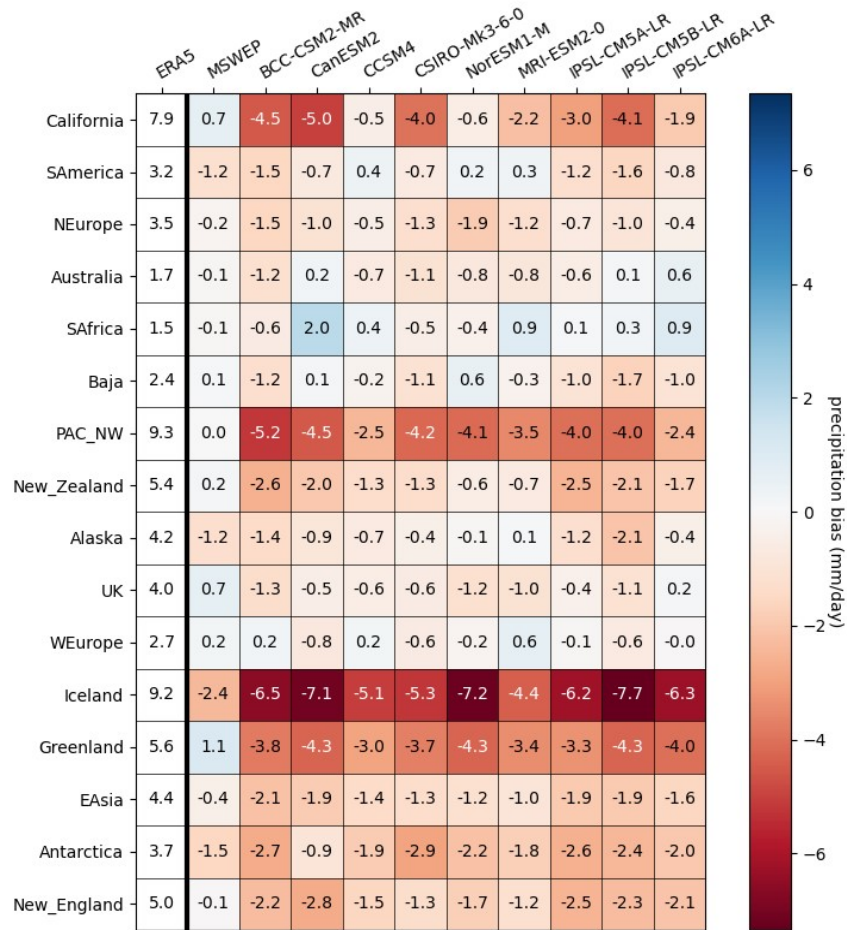
956
957
958



959
960
961
962
963
964
965
966

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

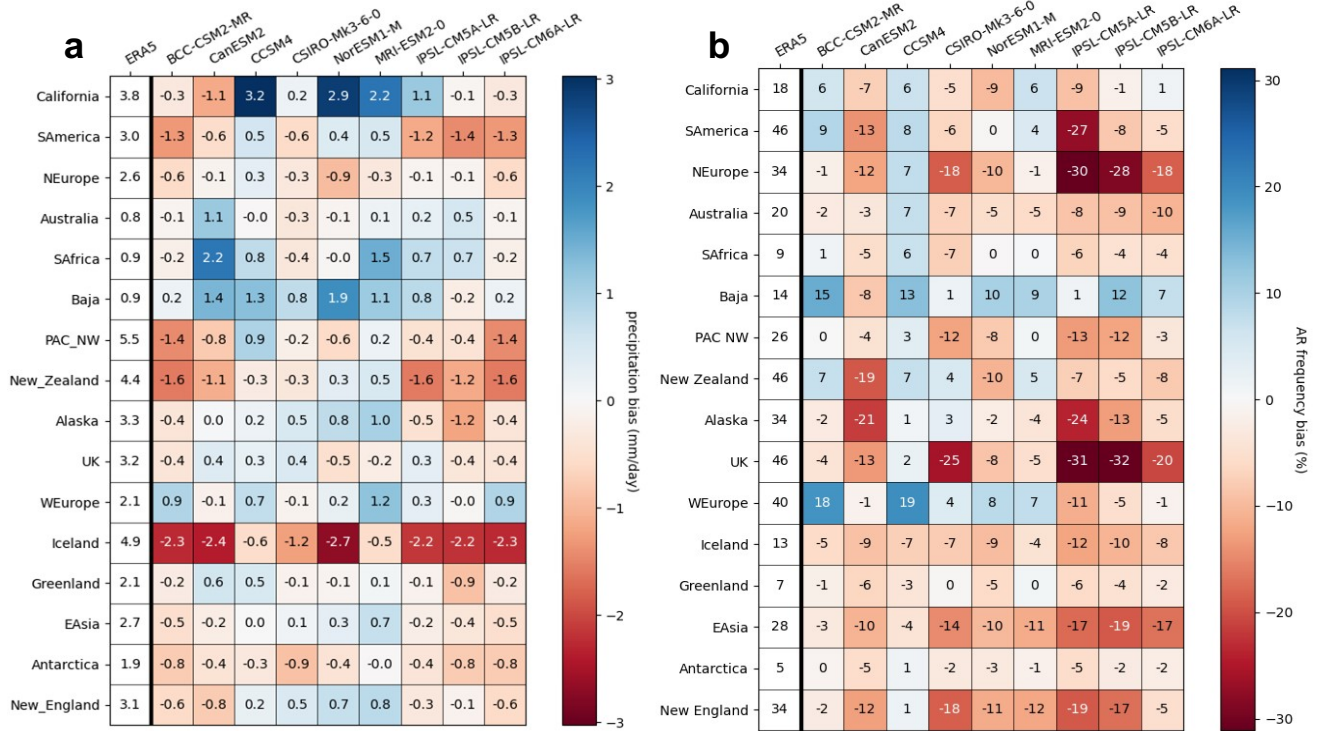
967
968



969
970
971

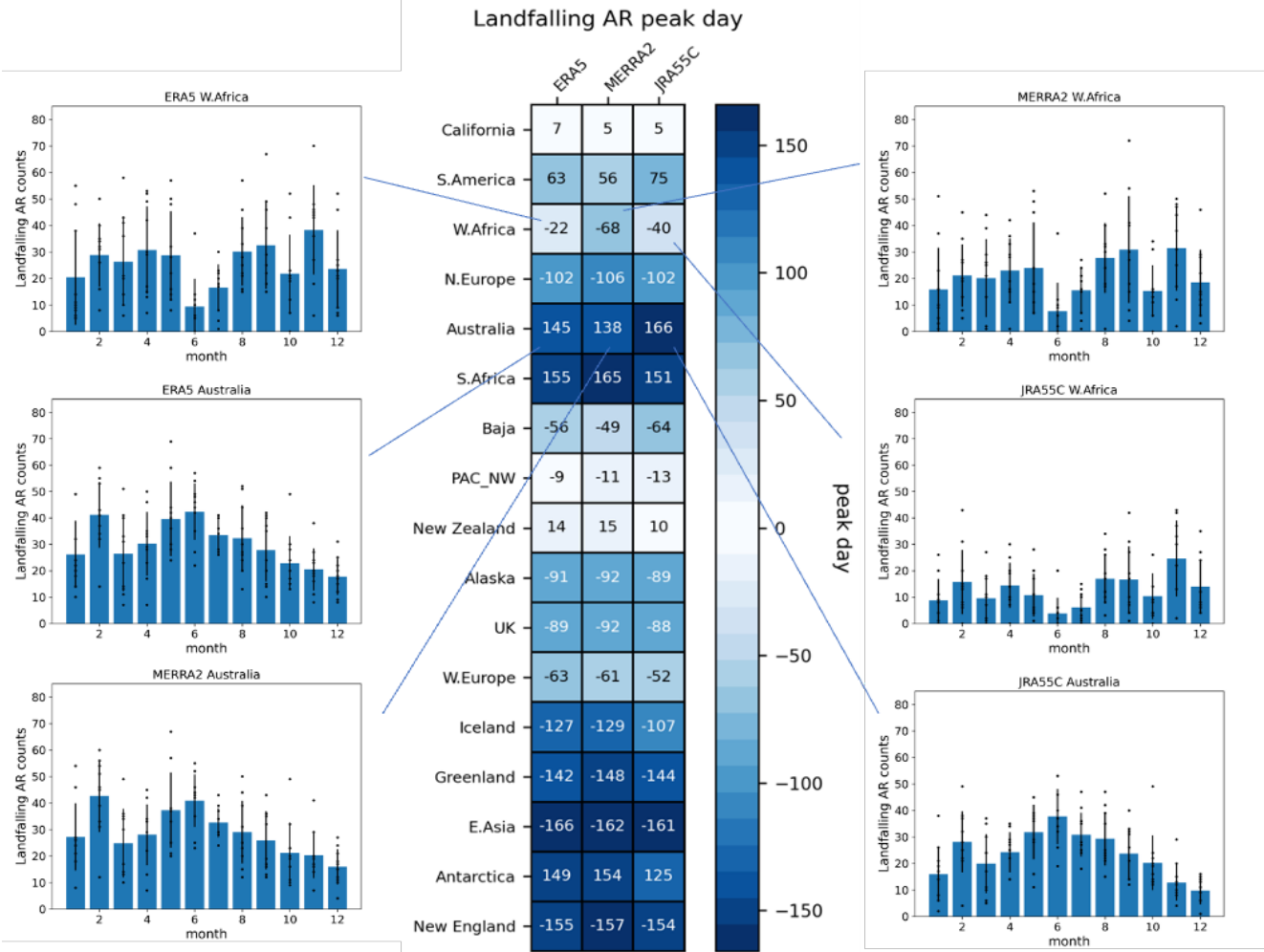
972 Fig. 10. Landfalling AR precipitation bias in climate models relative to ERA5
973 (the first column). The MSWEP data is also included in the second column
974 as an additional reference data, showed as the difference between ERA5.

975
976



977
978
979
980
981

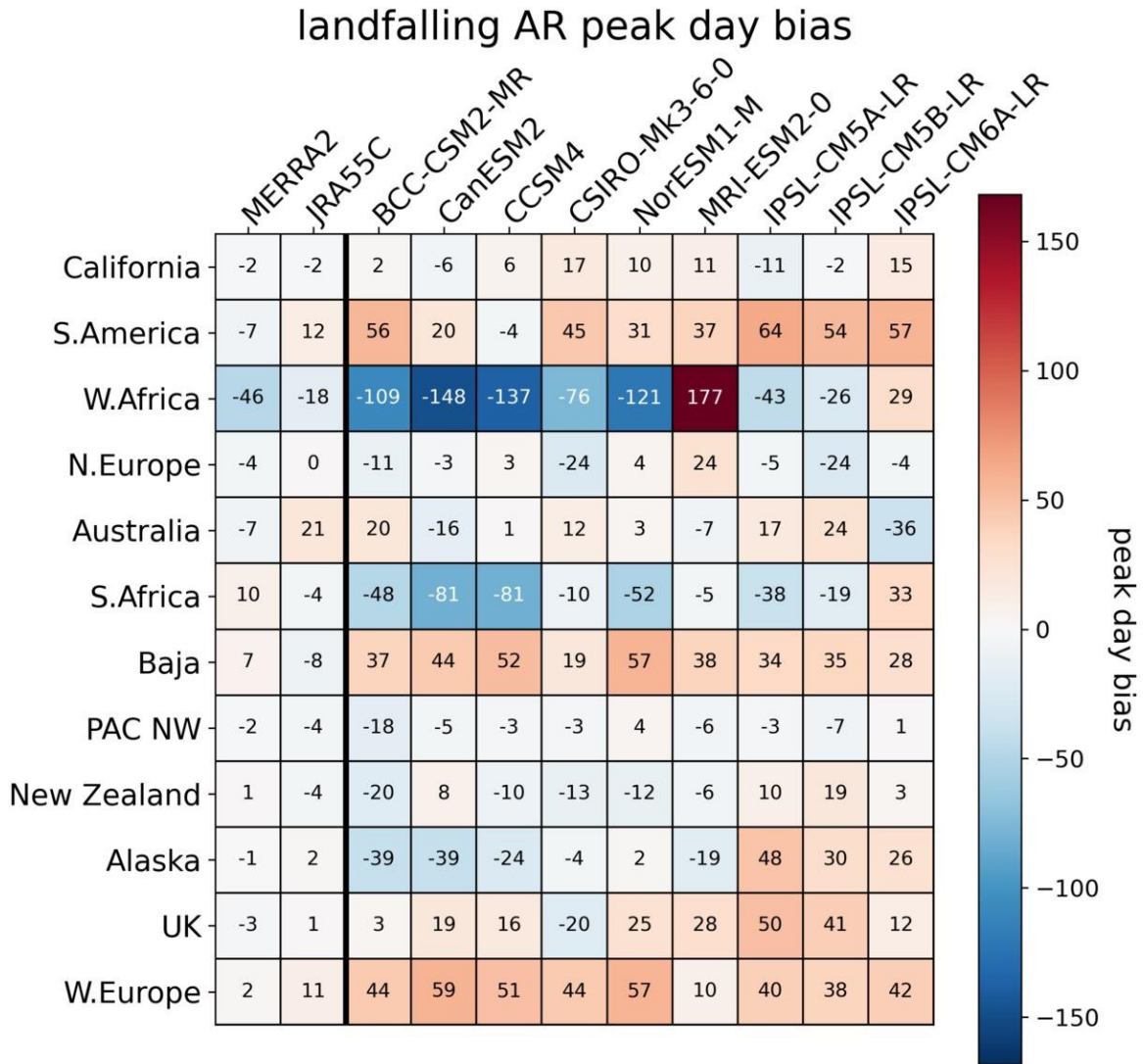
Fig. 11. (a) Total precipitation bias and (b) landfalling AR frequency bias



983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997

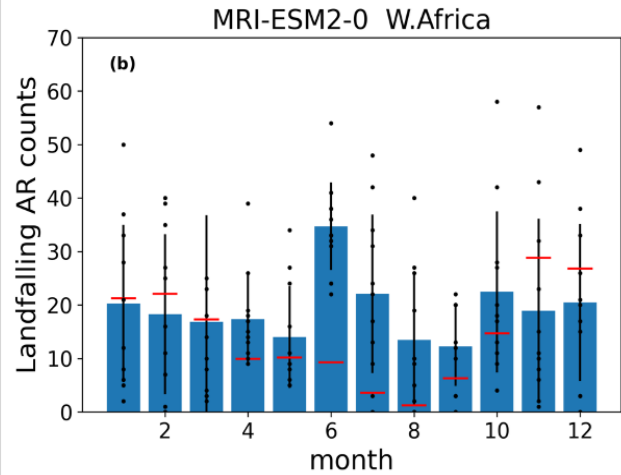
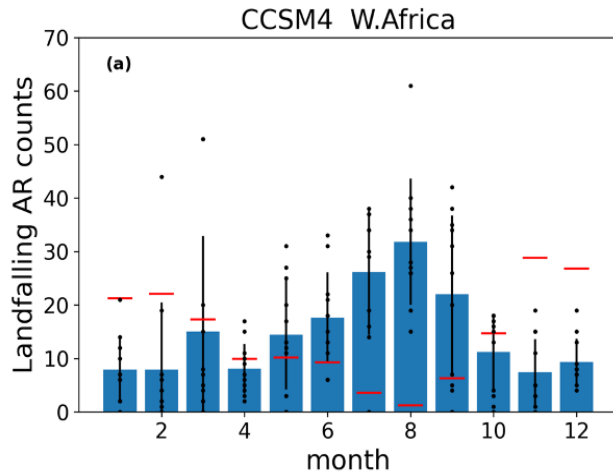
Fig. 12 (a) Landfalling AR peak day in ERA5, MERRA2, and JRA55C reanalysis. (b-g) show examples of probability distribution. Height of the blue bars indicate the time mean counts. Black dots represent peak day for each individual year, and vertical bars are the standard deviation range in the 10-year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

998
999



1000
1001
1002
1003
1004
1005
1006
1007
1008
1009

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



1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018

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.