



1 **Bergen Metrics: composite error metrics for assessing performance of climate models**
2 **using EURO-CORDEX simulations**

3 Alok K. Samantaray^{1,2}, Priscilla A. Mooney^{1,2}, Carla A. Vivacqua³

4 ¹Norwegian Research Centre (Norce), Norway

5 ²Bjerknes Centre for Climate Research, Norway

6 ³Federal University at Rio Grande do Norte, Brazil

7
8 Corresponding author:
9 Alok Kumar Samantaray
10 Jahnebakken 5, 5007 Bergen, Norway
11 Email: asam@norce-research.no

12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50



51

52 **Abstract**

53 Error metrics are useful for evaluating model performance and have been used extensively in
54 climate change studies. Despite the abundance of error metrics in the literature, most studies
55 use only one or two metrics. Since each metric evaluates a specific aspect of the relationship
56 between the reference data and model data, restricting the comparison to just one or two metrics
57 limits the range of insights derived from the analysis. This study proposes a new framework
58 and composite error metrics called Bergen Metrics to summarise the overall performance of
59 climate models and to ease interpretation of results from multiple error metrics. The framework
60 of Bergen Metrics are based on the p-norm, and the first norm is selected to evaluate the climate
61 models. The framework includes the application of a non-parametric clustering technique to
62 multiple error metrics to reduce the number of error metrics with minimum information loss.
63 An example of Bergen Metrics is provided through its application to the large ensemble of
64 regional climate simulations available from the EURO-CORDEX initiative. This study
65 calculates 38 different error metrics to assess the performance of 89 regional climate
66 simulations of precipitation and temperature over Europe. The non-parametric clustering
67 technique is applied to these 38 metrics to reduce the number of metrics to be used in Bergen
68 Metrics for 8 different sub-regions in Europe. These provide useful information about the
69 performance of the error metrics in different regions. Results show it is possible to observe
70 contradictory behaviour among error metrics when examining a single model. Therefore, the
71 study also underscores the significance of employing multiple error metrics depending on the
72 specific use case to achieve a thorough understanding of the model behaviour.

73

74

75

76

77

78

79

80

81

82

83

84



85 1. Introduction

86 Climate models are important tools for predicting and understanding climate change, and
87 climate processes (Kotlarski et al., 2014; IPCC, 2021a; IPCC, 2021b; Mooney et al., 2022). In
88 the context of climate studies, climate model evaluation is essential for identifying models that
89 poorly simulate the climate system, and for ranking of climate models (Randall et al., 2007;
90 Flato et al., 2013). The main purpose of climate model evaluation is twofold; firstly, to ensure
91 that the models are reproducing key aspects of the climate system and secondly to understand
92 the limitations of climate projections from the models. This ensures proper interpretation and
93 application of climate models and any climate projections produced by them. The performance
94 of climate models is quantified by different error metrics such as root mean square error, and
95 bias, which assess the agreement between the climate model data and reference data (e.g.,
96 gridded observational products, station data, reanalyses, or satellite observations).
97 Different error metrics are available in the literature, and each has a specific framework
98 according to its purpose (Rupp et al., 2013; Pachepsky et al., 2016; Baker & Taylor, 2016;
99 Collier et al., 2018; Jackson et al., 2019). For example, root mean square error compares the
100 amplitude difference between modelled and reference data, while the correlation coefficient
101 compares the phase difference between modelled and reference data. Depending on the specific
102 error, the error metrics can be categorised into different classes; the most popular classes are
103 accuracy, precision, and association. Accuracy measures the degree of similarity between
104 climate model data and reference data. An extremely high accuracy indicates that the model
105 has less error magnitude of any type and testing the model with other error metrics adds little
106 value (Liemohn et al., 2021). However, if a model has moderate to low accuracy, testing the
107 model with other metrics can reveal other similarities and dissimilarities between model data
108 and reference data. Root mean square error and mean square error are the most used accuracy
109 metrics to evaluate climate models (Watt-Meyer et al., 2021; Wehner et al., 2021; He et al.,
110 2021), even though the metrics cannot reveal whether the model is under or over-predicting
111 the observations. Precision metrics quantify the degree of similarity in the spread of the data.
112 A robust and commonly used metric for assessing the precision of model data is the ratio or
113 difference of standard deviation between modelled data and reference data (van Noije et al.,
114 2021; Wood et al., 2021; Wehner et al., 2021). Finally, association metrics measure the degree
115 of the phase difference between modelled data and observed data. Phase difference is important
116 in climate studies as it affects the initiation and termination time of a season of climate
117 variables. One metric that is extensively used to measure the association is the correlation
118 coefficient (Richter et al., 2022; Bellomo et al., 2021; Yang et al., 2021). Liemohn et al. (2021)



119 has described various other major categories of metrics and they suggest that assessment of
120 models should not be restricted to one or two error metrics. Interested readers can follow the
121 citations to read in detail about the discussed metrics.

122 There are several composite error metrics that use the modified framework of other metrics to
123 compute the error magnitude. A widely used example of this is the Taylor diagram (Taylor,
124 2001), which incorporates correlation, root mean square deviation and ratio of standard
125 deviation. A distinguishing feature of the Taylor Diagram is its ability to graphically evaluate
126 the model performance. Another popular example is the Nash-Sutcliffe Efficiency (NSE; Nash
127 & Sutcliffe, 1970) which is a normalised form of the mean squared error to evaluate and predict
128 the model streamflow data. Later, it was observed that NSE can be decomposed into three
129 components which are the functions of correlation, bias and standard deviation (Murphy, 1988;
130 Węglarczyk, 1998). Other similar scores include the Kling-Gupta (K-G) efficiency (Gupta et
131 al., 2009) which is a function of three components: ratio of model mean to observed mean, the
132 ratio of model standard deviation to observed standard deviation and correlation coefficient.

133 The study of Gupta et al. (2009) argued the NSE, which has a bias component normalised by
134 the standard deviation of the reference data, will have a low weight on the bias component if
135 the reference data has high variability. The modified Kling-Gupta efficiency developed by
136 Kling et al. (2012) involves the ratio of covariance instead of the ratio of standard deviation.

137 Both K-G efficiency and modified K-G efficiency use Euclidean distance as a basis to calculate
138 the error magnitude of the model and the study argued that instead of finding a corrected NSE
139 criterion, the whole problem can be viewed from the multi-objective perspective where the
140 three error components can be used as separate criteria to be optimised. It identifies the best
141 models by calculating the Euclidean distance from the ideal point and then finding the model
142 with the shortest distance. The ideal value of an error metric is obtained when the model exactly
143 simulates the observed data. The Euclidean distance is also used by Hu et al. (2019) to develop
144 the DISO metric that incorporates correlation coefficient, absolute error and root mean squared
145 error. The study of Hu et al. (2019) also argues that accuracy (root mean square error), bias
146 (absolute error) and association (correlation coefficient) are the three major error classes based
147 on which a model should be assessed and evaluating a model using a single error metric may
148 lead to ill-informed results. The study pointed out a few limitations of the Taylor diagram such
149 as quantification of error magnitude and low sensitivity to small error differences by the
150 diagram. In a comparative study, Kalmár et al. (2021) found no substantial difference between



151 DISO index and the Taylor diagram. However, based on quantification of error magnitude,
152 DISO index can be helpful.

153 The Euclidean distance framework has been increasingly used in different fields as an error
154 function or metric for many applications such as evaluation of models, parameter
155 optimization and classification problems. Euclidean distance is basically the second norm of a
156 vector. Equation 1 is the generalised form of p-norm in a n-dimensional vector space, where
157 x_i is the vector. When p is 2, it becomes the Euclidean norm. If the vector (x_i) is the difference
158 between the observed data (u_i) and model data (v_i) i.e. $x_i = u_i - v_i$, then d is called the
159 Euclidean distance metric. i represent the time series data. Root mean squared error and mean
160 squared error are different variants of Euclidian distance metric. If the vector is the difference
161 between error metrics (correlation coefficient [u_1], absolute error [u_2] and root mean squared
162 error [u_3]) and their ideal values ($v_{1:3}$), then d is called the DISO index. A disadvantage of the
163 Euclidean distance is that it suffers the curse of dimensionality (Mirkes et al., 2020; Weber et
164 al., 1998) i.e. Euclidean distance as a dissimilarity index becomes less efficient as dimension
165 increases. In this study, we assess the effect of the norm order on the overall error. We use
166 different measures such as the contribution of outliers to the overall error, the difference
167 between the maximum and minimum distances, and the average distances to compare different
168 norms.

$$169 \quad d_n(u, v) = (\sum_{i=1}^n |x_i(u_i, v_i)|^p)^{1/p} \quad (1)$$

170 This study has the following objectives:

- 171 i) Evaluation of 89 CMIP5 driven regional climate simulations from the Euro-
172 CORDEX initiative using 38 error metrics;
- 173 ii) Clustering of error metrics to assess their performance;
- 174 iii) Assessment and recommendation of different p-norms based on their performance;
- 175 iv) Formulation of a composite metric using the optimal norm.

176 **2. Data and Study area**

177 We focus on Europe due to the widespread availability of a large ensemble of high resolution
178 (0.11°) regional climate simulations. In this study, we use 89 regional climate model (RCM)
179 simulations from Euro-CORDEX to study the behaviour of different error metrics. The Euro-
180 CORDEX dataset provides both precipitation and temperature data at 0.11° grid resolution.
181 The monthly data from 1975 to 2005, which is available in all the RCM simulations, have been
182 used to calculate the index. Supplementary Table S1 provides an overview of the global climate



183 models (GCMs) downscaled by the different RCMs. Supplementary Table S2 provides an
184 overview of the RCMs and assigns a number (Column 1) to each RCM which is used to identify
185 RCMs in plots that have limited space for labels.

186 For reference data, both precipitation and temperature data are obtained from E-OBS dataset.
187 The reference data has a 0.25° grid spacing. To compare the model data with the reference
188 data, all the data needs to be on a common grid. In this study, we remapped the RCM data onto
189 the coarser 0.25° grid of E-OBS.

190 The study uses the eight sub-regions of Europe defined by Christensen & Christensen (2007)
191 – British Isles, Iberian Peninsula, France, Mid-Europe, Scandinavia, Alps, Mediterranean, and
192 Eastern Europe - to conduct analysis in more homogeneous areas.

193 **3. Methodology**

194 This section outlines the framework for clustering error metrics and provides a brief overview
195 of their characteristics. Additionally, the section describes the proposed metric's framework.

196 **3.1 Error metrics**

197
198 Error metrics are commonly used in climate change studies to measure the differences between
199 modelled and reference data in time series. As the number of climate models has increased, the
200 study of error metrics has become increasingly important. There are several error metrics
201 available to evaluate the performance of climate models (Jackson et al., 2019), and the selection
202 of an appropriate metric remains a topic of debate in the literature. For instance, Willmott &
203 Matsuura (2005) advocate for mean absolute error (MAE) over root mean squared error
204 (RMSE), as the latter is not an effective indicator of average model performance. In contrast,
205 Chai & Draxler (2014) contend that RMSE is superior to MAE when errors follow a Gaussian
206 distribution. To gain insight into the performance of error metrics, we have analysed Euro-
207 CORDEX precipitation data and examined the differences in ranking of 89 GCM-driven
208 regional climate simulations using 38 error metrics (Jackson et al., 2019). The list of error
209 metrics is provided in Table S3. All 89 models are ranked based on their performance using
210 the 38 error metrics. The average ($r_{M,mean}$; Equation 2) and maximum ($r_{M,max}$; Equation 3)
211 rank differences are then calculated at each grid point. The former is the mean of all the
212 pairwise rank differences, while the latter is the maximum of all the pairwise rank differences.
213 These calculations allow us to understand the performance of different error metrics and the
214 extent of the disparity in ranking of the climate models.

215



216 **Table 1: Example of ranking order**

Number	Climate model	Ranking order (RO) by i th error metric (E_i)	Ranking order (RO) by k th error metric (E_k)
1	M1	3	2
2	M2	1	3
3	M3	2	1

217

$$218 \quad r_{M,mean} = \mu_g(R_{M,k} - R_{M,i}) \quad (2)$$

$$219 \quad r_{M,max} = \max_g(R_{M,k} - R_{M,i}) \quad (3)$$

220 $R_{M,k}$ and $R_{M,i}$ are the rank assigned to model M by the k th and i th error metric, respectively.

221 We have provided Table 1 as an example for better understanding of the notations. If there are

222 three climate models (M1, M2 and M3) as shown in Table 1, all the models have been assigned

223 to a number (first column) and the order must not change throughout the study. $R_{M,k}$ and $R_{M,i}$

224 for model M1 are 2 and 3, respectively. k varies from 1 to N_E-1 and i varies from $k+1$ to N_E ,

225 where N_E is the total number of error metrics. The difference in ranking is calculated for all

226 possible combinations of error metrics. $\mu_g()$ and $\max_g()$ are the mean and maximum operator,

227 respectively, which is applied across all the grid points ($g:1,2,\dots,gd$). gd is the total number of

228 grid points which is 11370 in this study. Figure 1 demonstrates that different error metrics used

229 to assess climate models result in significantly different ranking orders. The average of $r_{M,mean}$

230 across all the grid point varies from 16 to 26 whereas the average of $r_{M,max}$ varies from 40 to

231 70. The results indicate significant differences in the ranking of the climate models by different

232 error metrics. The disparity in ranking order may be due to the distinctive error targeted by

233 each metrics as discussed in the introduction section.

234 This study assumes that all the errors are important and that it may be necessary to evaluate

235 model performance using multiple metrics. To achieve independence among the metrics, the

236 study has attempted to cluster the error metrics based on model performance. This classification

237 would enable different clusters to have unique characteristics, and metrics within the same

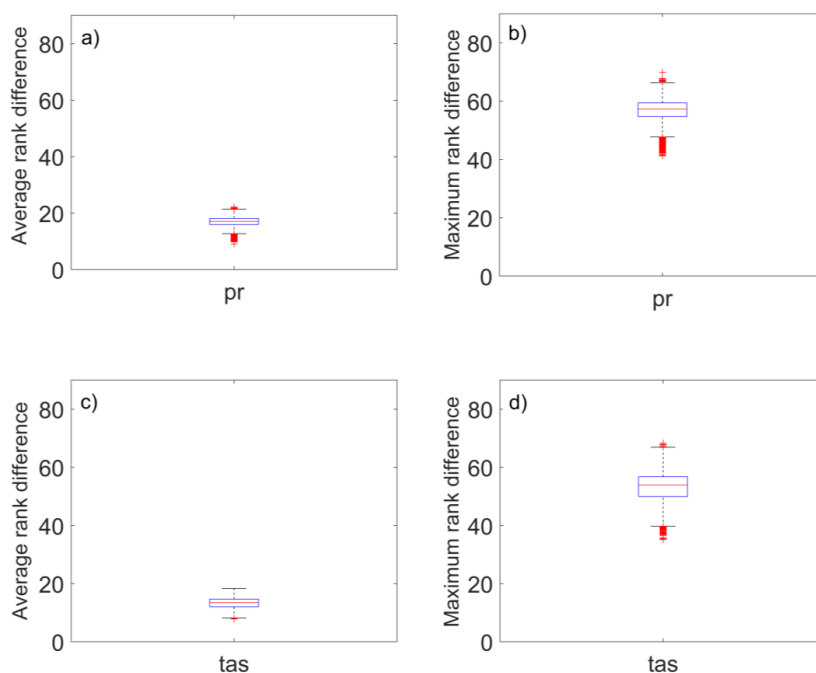
238 cluster would produce similar results, whereas those from different clusters would yield

239 different ranking orders. In summary, the study proposes that using multiple error metrics and

240 clustering them based on performance could improve the understanding and

241 comprehensiveness of climate model analysis.

242



243

244 **Figure 1:** Box plot of average rank difference (first column [a, c]) and maximum rank
245 difference (second column; [b, d]) for precipitation (Pr; first row [a, b]) and temperature (T;
246 second row [c, d]) over all the grid points in European region

247 3.2 Clustering of error metrics

248 The aim of clustering error metrics is to group a set of metrics based on their similarities such
249 that the metrics within the same cluster generate similar rankings of climate models compared
250 to those in different clusters. This study clusters the error metrics using a non-parametric
251 clustering approach inspired by the Chinese restaurant process (CRP; Pitman, 1995). This
252 approach was chosen based on its performance compared to the k-means clustering approach
253 (see Text S1) and its simpler framework. The algorithm follows two fundamental principles:
254 (i) the first error metric (E_1) forms the first cluster (C_1), and (ii) the i th error metric (E_i) is
255 assigned to a cluster which has the maximum of all the mean absolute error (u_j) values greater
256 than a particular threshold value (th). The clustering algorithm is presented in Fig. 2.

257 Similar to the rank difference explained in the previous section, the MAE (RO_i, RO_k) between
258 the ranking order produced by two error metrics is computed. RO is the ranking order and it
259 can be calculated by assigning the climate models to a number. For example, the ranking order
260 (RO_i) by i th error metric and the ranking order (RO_k) by k th error metric are [3, 1, 2] and [2,
261 3, 1], respectively in Table 1. The MAE values are calculated for all possible combinations of



262 error metrics in a particular cluster and the maximum of the MAE values is used to compare it
263 to the threshold value. The exercise is repeated for all the clusters (N_C) available at that time.
264 The number of clusters (N_C) and the number of error metrics in each cluster (N_{CE}) are updated
265 for each iteration (i) and if the criteria is not satisfied, then a new cluster is formed using that
266 error metric. The whole exercise is repeated till all the error metrics (N_E) gets assigned to a
267 cluster.

```
E1 ∈ C1           First error metric belongs to the first cluster
For i = 2:NE do      For all the error metrics
  For j < NC do      For all the clusters
    For k < NCE do   For all the error metrics in Cj
      Uj,k = MAE(ROi,ROk)
    uj = max(Uj,k)
  If uj < th
    Ei ∈ Cj
  else
    Ei ∈ CNC+1
```

268

269 **Figure 2:** Algorithm of the non-parametric clustering for classifying the error metrics

270 The threshold value is defined as qth percentile of a column matrix D where D is the collection
271 of MAE values for all possible combinations of error metrics at all the grid points in a region.
272 In this study, q has been assigned the value of 10 and the sensitivity of q is discussed in the
273 results section.

274 3.3 Proposed metric- The Bergen Metrics

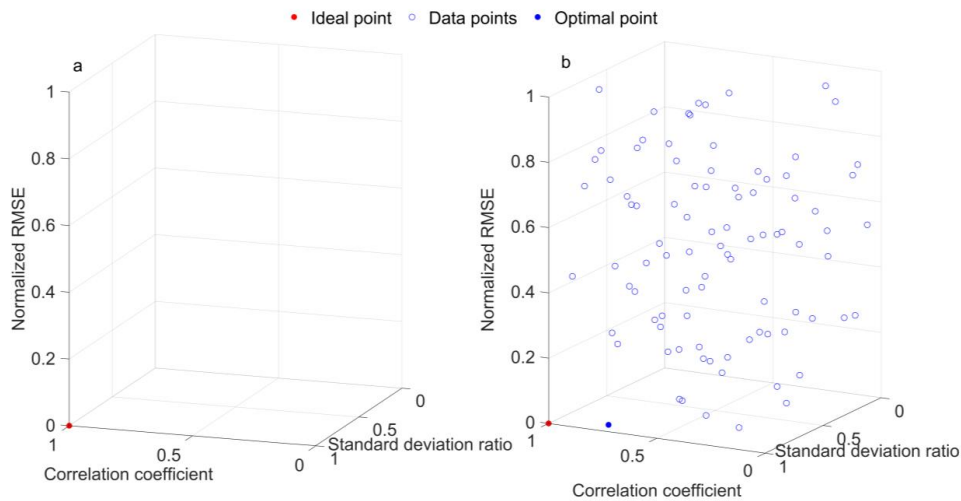
275 The clustering of error metrics guarantees that metrics in different groups produce distinct
276 ranking orders, implying that each group targets different errors. One of the objectives of this
277 study is to integrate different errors and create a composite error to obtain a single value. One
278 potential solution is to use the Euclidean distance approach with different error metrics as
279 different dimensions in the Euclidean space. To illustrate this, we employed three widely used
280 error metrics: Normalized Root Mean Square Error (RMSE), Standard Deviation ratio (SD)
281 and correlation coefficient. In the Euclidean space, an ideal model that predicts the climate
282 variable as accurately as the observed data would have values of 1, 1, and 0 for correlation
283 coefficient, Standard Deviation ratio, and normalized RMSE, respectively. The coordinates of



284 an ideal model in the Euclidean space would be (1, 1, 0), as represented by the red point in Fig.
 285 3a. Since different models have unique coordinates based on the three metrics, these
 286 coordinates serve as possible solutions to determine the best model. If a decision is required,
 287 one approach could be to calculate the Euclidean distance from the ideal point to all points and
 288 select the point with the shortest distance (Equation 4). This equation can be simplified to
 289 Equation 5. The model that is closest to the ideal point, indicated by the optimal point in Fig.3b,
 290 can be considered as the best model.

291
$$ED\ Metric = \sqrt{(1 - Correlation\ coefficient)^2 + (1 - Standard\ deviation\ ratio)^2 + (0 - RMSE)^2} \quad (4)$$

292
$$ED\ Metric = \sqrt{(1 - Correlation\ coefficient)^2 + (1 - Standard\ deviation\ ratio)^2 + (RMSE)^2} \quad (5)$$



293
 294 **Figure 3:** Example for three-dimensional (a) ideal point and (b) the solution space of
 295 correlation coefficient (x-axis), standard deviation (y-axis) and normalized RMSE (z-axis)
 296

297 The Euclidian distance has several benefits that make it a popular metric, primarily its
 298 simplistic framework. However, it also has some drawbacks. The Euclidian distance, also
 299 known as L2 norm, is less effective in higher dimensional spaces, which can lead to instability
 300 when additional error metrics are added (Weber et al., 1998; Aggarwal et al., 2001). To mitigate
 301 this issue, recent research has focused on the use of L1 norms, such as relative mean absolute
 302 error and mean absolute scaled error, which have become more popular than L2 norms like
 303 mean squared error. This approach reduces the impact of outliers in the data (Armstrong &



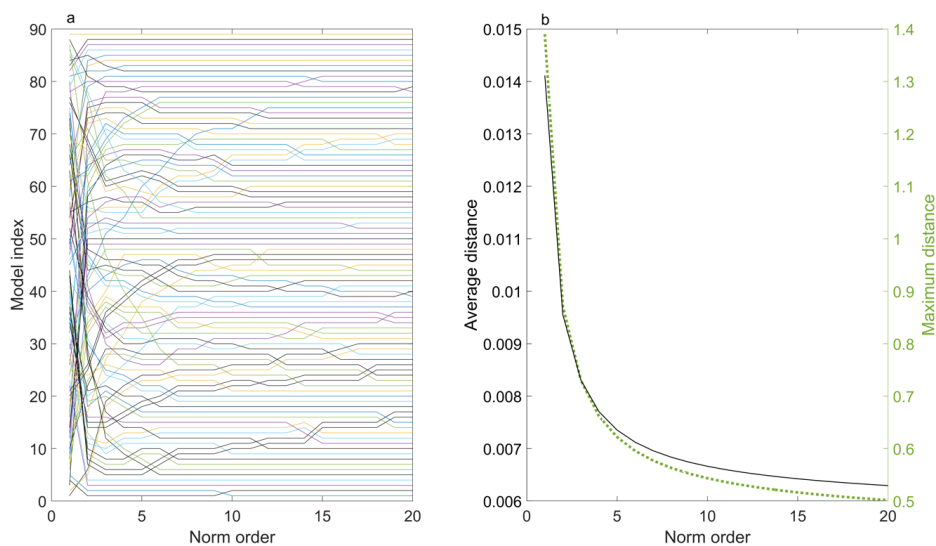
304 Collopy, 1992; Hyndman and Koehler, 2006). Reich et al. (2016) found that relative MAE,
305 based on an L1 norm, is advantageous in assessing prediction models. This study proposes the
306 following new metrics called the Bergen Metrics (BM) which is a generalised p-norm
307 framework to evaluate climate models. Equation 6 presents the generalised form of the metric.
308 It is important to note that equation 6 serves as an illustration of Bergen metrics, and users
309 have the flexibility to include or remove metrics according to their preference.

$$310 \text{ Bergen Metric} = \sqrt[p]{\begin{matrix} (1 - \text{Correlation coefficient})^p \\ + (1 - \text{Standard deviation ratio})^p \\ + (0 - \text{RMSE})^p \end{matrix}} \quad (6)$$

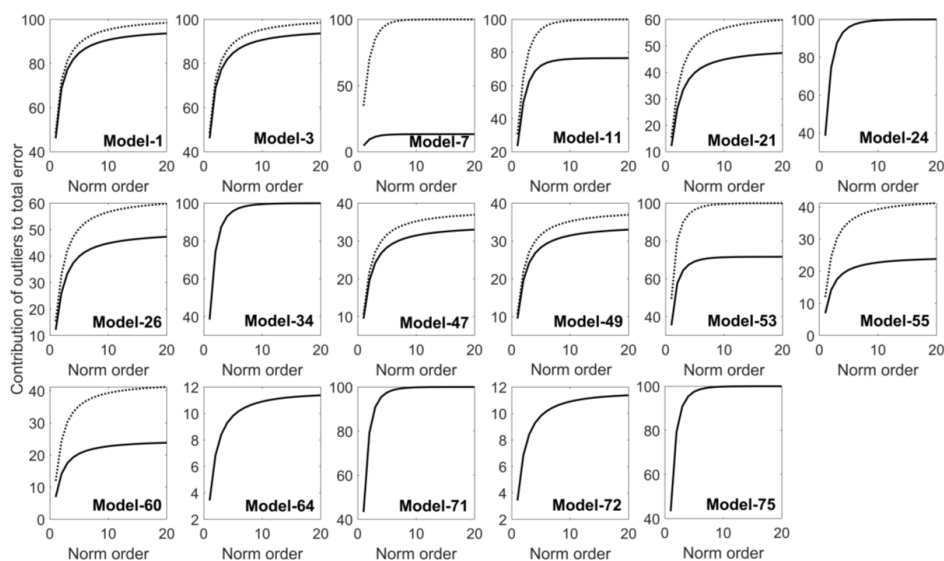
311 A case study has been conducted to understand the impact of different p norms on the ranking
312 order of climate models. For this, five error metrics - RMSE, bias, correlation coefficient,
313 standard deviation ratio, and mean ratio - have been considered (Equation 7) and the error
314 metrics are normalised using model data. The study includes 89 RCM simulations for
315 precipitation, and Fig. 4a shows the ranking of these models for different p norms. The lines
316 corresponding to each model give information about the model's ranking in different norms.
317 The results demonstrate that climate models are highly sensitive to p norms. Significant change
318 in ranking order is observed for the first four norms. Fig. 5 shows the percentage contribution
319 of outliers to the total error magnitude for models that have outliers. Median absolute deviation
320 technique (MAD) is used to identify outliers among the error metrics. Some of the models
321 have only one outlier (plots with a single solid line in Fig. 5) and other models have two outliers
322 (plots with both solid and dotted lines in Fig. 5). The percentage contribution of outliers
323 increases as the p norm increases, consistent with previous literature (Armstrong and Collopy,
324 1992; Hyndman and Koehler, 2006). The study has used two parameters to indicate the
325 capability of each norm to differentiate between climate models - mean pairwise difference of
326 the BM and the difference between the maximum and minimum values of the BM. Figure 4b
327 shows that both parameters decrease as the p norm increases, indicating less differentiability.
328 The results suggest that the first norm (p=1) is the optimal norm to use as a metric in this study
329 and will be utilized in the following analyses.

$$330 \text{ Bergen Metric (BM)} = \sqrt[p]{\begin{matrix} (0 - \text{RMSE})^p + (0 - \text{Bias})^p \\ + (1 - \text{Standard deviation})^p \\ + (1 - \text{Correlation coefficient})^p + (1 - \text{Mean ratio})^p \end{matrix}} \quad (7)$$

331



332
333 **Figure 4:** a) The change in the ranking of the climate models with different norm order (p) b)
334 the change in the difference between the maximum and minimum distances and the average
335 distances with different norm order



336
337 **Figure 5:** The percentage contribution of outliers to the total error magnitude as a function of
338 norm order. The colours represent different outliers.

339



340 4. Results

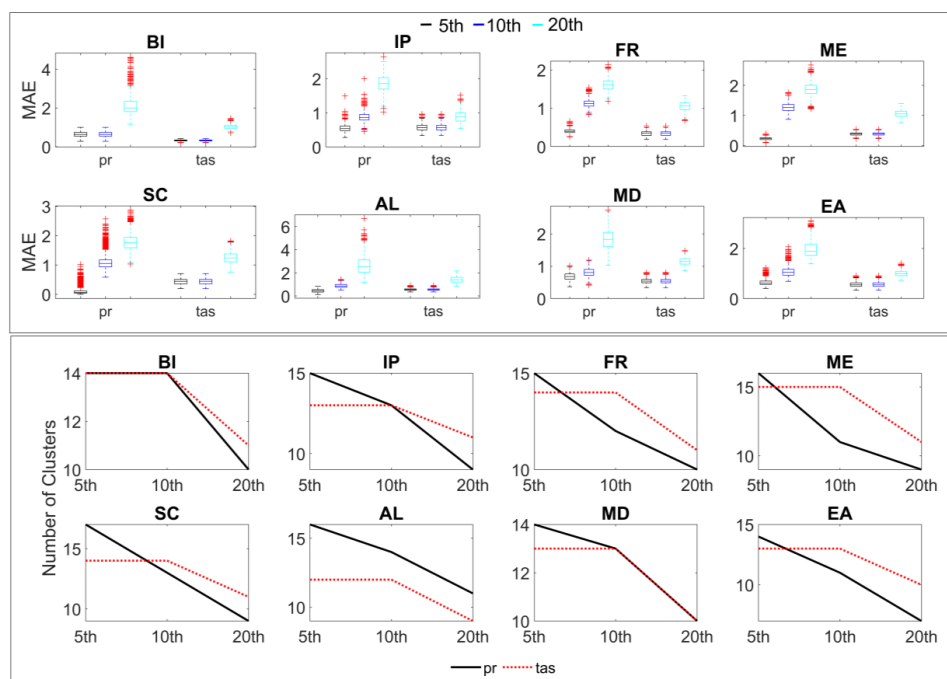
341 4.1 Regional clustering of error metrics

342 The study considers 38 error metrics (Table S3) which can take both positive and negative
343 values as input. Similar to the models, the error metrics have been assigned a number (column
344 1; Table S3) and the error metrics have been labelled as those numbers in some figures.

345 The clustering technique described in the methodology section can be applied to individual
346 grid points, but for the sake of simplicity, we use a single cluster for all grid points within each
347 of these regions defined by Christensen & Christensen (2007). The methodology is modified
348 slightly to enable regional clustering. At a grid point scale, the maximum value of mean
349 absolute error (u_j) is used as a proxy for that specific error metric at a grid point. For regional
350 clustering, the maximum MAE values are computed for all grid points within the region, and
351 the average of those values is used as a proxy for that region and error metric. This value is
352 then compared with a threshold to determine whether the error metric belongs to a certain
353 cluster or it should be assigned to a new cluster. The clustering algorithm is executed for
354 multiple thresholds.

355 The 5th, 10th, and 20th percentiles are selected as potential thresholds to cluster the error
356 metrics. However, users can select any number of thresholds for the sensitivity analysis. The
357 clustering algorithm is allowed to run for all the thresholds to determine the optimal threshold.
358 The efficiency of each cluster for a given threshold is represented by the mean of MAE over
359 all the clusters. Another criterion used to determine the threshold is the number of clusters
360 corresponding to each threshold. An increase in the percentile (q) is expected to increase the
361 MAE as the magnitude of threshold increases. Similarly, the number of clusters are expected
362 to decrease as q increases as it can allow more error metrics into a cluster due to higher
363 threshold magnitude. From Fig. 6, we conclude that the results are according to our
364 expectations. It is found that increasing the percentile resulted in an increase in MAE and a
365 decrease in the number of clusters. The 10th percentile is selected as the threshold to cluster
366 the error metrics for both temperature and precipitation, as it has a smaller number of clusters
367 compared to 5th percentile and less MAE compared to 20th percentile. The

368



369

370 **Figure 6:** The variation in MAE (first box) and number of clusters (second box) corresponding
371 to 5th, 10th and 20th percentile for precipitation (pr) and temperature (tas) for all the eight regions

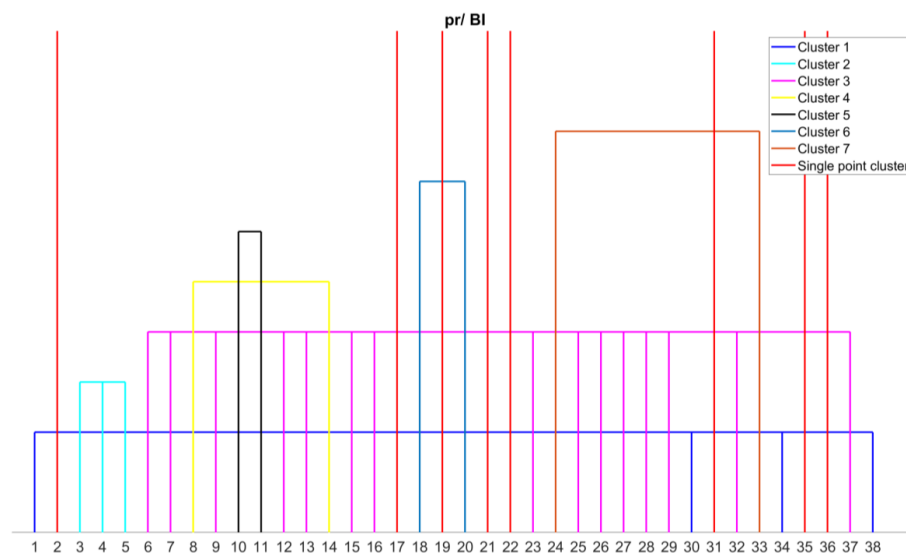
372 4.2 Results of clustering

373 4.2.1 Precipitation

374 For the British Isles region, the classification of 38 error metrics resulted in 15 clusters, with 8
375 error metrics being single point clusters due to their unique behaviour (Fig. 7). These 8 metrics
376 are d [2], (MB) R [17], MdE [19], MEE [21], MV [22], r2 [31], SGA [35], and R(Spearman)
377 [36]. The threshold for precipitation data is 6.35, indicating that all 8 error metrics produced
378 MAE values greater than 6.35 compared to the remaining 30 error metrics. RMSE [32] and its
379 variants such as normalized RMSE by IQR [25], mean [26] and range [27] are assigned to the
380 same cluster, as ED [7], IRMSE [9], MAE [13], MAPD [15], MASE [16], and MSE [23]. The
381 reason could be the L-norm framework which is used by most of the error metrics in this cluster.
382 D1 [3], d1 [4], and d(Mod.) [5] which share a similar framework, are also assigned to a single
383 cluster. Error metrics that evaluate the phase difference between observed and modelled data,
384 including ACC [1], R (Pearson) [30], SC [34], and M [38], are assigned to a single cluster.
385 H10(MAHE) [8] and MALE [14] share the same cluster as both metrics consider the difference
386 of logarithmic of the model and observed data to compute the error. Similarly, MdAE [18] and



387 MdSE [20] are assigned to a single cluster, as both metrics use the median of the difference
388 between observed and modelled data. However, MdE [19] is assigned to a different cluster as
389 it only considers the difference between observed and modelled data without bringing them to
390 the positive domain. NED [24] and SA [33] are found to be in the same cluster, as both metrics
391 are linearly associated while evaluating the model, even though their underlying frameworks
392 are somewhat different. Although ED [7] and NED [24] follow the L2 norm, they are not
393 assigned to the same cluster. This can be attributed to the normalisation of observed and
394 modelled data by their respective means in NED, as the statistical parameters such as mean is
395 sensitive to outliers, which can result in changes in ranking order.
396



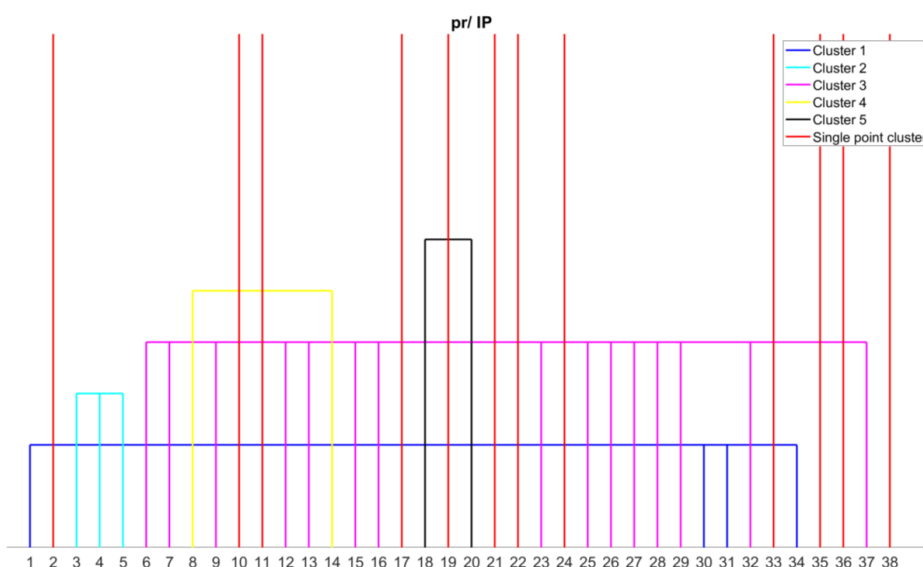
397

398 **Figure 7:** Clustering of error metrics using precipitation (pr) data for British Isles (BI) region.
399 Each error metric can be identified by the number using Table S3.

400 The Iberian Peninsula region is found to have 17 clusters, with 12 of them being single point
401 clusters (Fig. 8). Seven of the eight error metrics that are single point clusters in British Isles
402 are also single point clusters in Iberian Peninsula, except for r2 [31]. Five other error metrics:
403 NED [24], KGE (2009) [10], KGE (2012) [11], SA [33], and M [38] are also single point
404 clusters in Iberian Peninsula region. In British Isles, KGE (2009) [10] and KGE (2012) [11]
405 are assigned to the same cluster. The KGE (2012) is different from KGE (2009) since it used
406 the ratio of coefficient of variation between modelled and observed data instead of the ratio of
407 standard deviation to avoid the cross-correlation between bias and variability ratio. The



408 coefficient of variation is the ratio between the standard deviation and the mean of the data,
409 which represents the extent of variability with respect to the mean of the data. A biased dataset
410 can produce a significant change in the relative standard deviation, i.e., the coefficient of
411 variation. That is a possible reason why both the metrics are in different clusters. r^2 is assigned
412 to the correlation metrics cluster in this region. The remaining clusters are almost identical to
413 the clusters obtained for the British Isles region.



414

415 **Figure 8:** Clustering of error metrics using precipitation (pr) data for Iberian Peninsula (IP)
416 region. Each error metric can be identified by the number using Table S3.

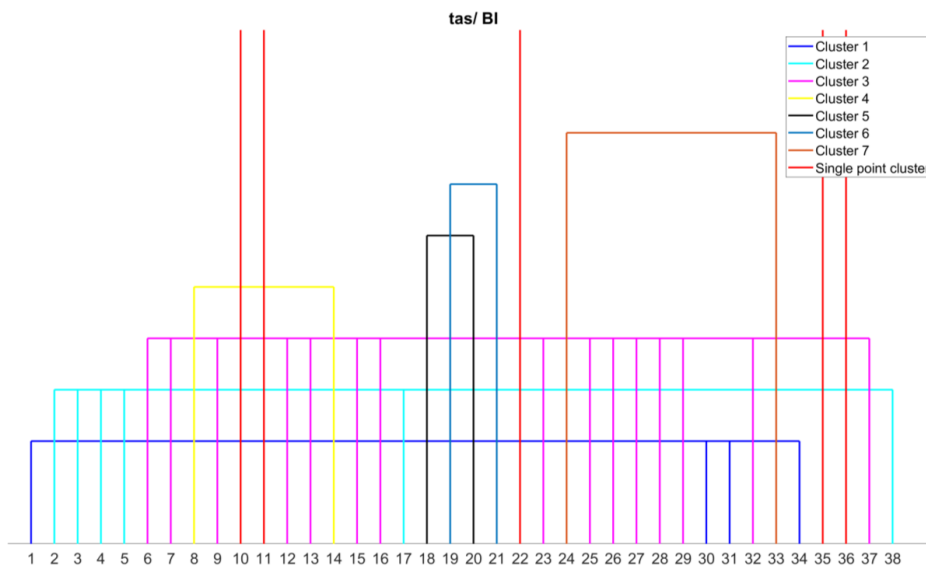
417 As the results for the other 6 regions are similar to either the British Isles or the Iberian
418 Peninsula, we simply summarise their results here and refer the reader to the supplementary
419 material for further information. France (Fig. S2), Mid-Europe (Fig. S3), Scandinavia (Fig.
420 S4), Alps (Fig. S5), Mediterranean (Fig. S6) and Eastern Europe (Fig. S7) exhibit 15, 15, 16,
421 16, 17, and 14 clusters, respectively, with 8, 8, 10, 10, 12, and 6 single point clusters. France
422 and Mid-Europe have the same clusters as the British Isles, and the Mediterranean has the same
423 clusters as Iberian Peninsula. Scandinavia has clusters similar to British Isles, except that M
424 [38] is a single point cluster and r^2 [31] has been assigned to the correlation metrics cluster in
425 Scandinavia. The Alps also has clusters similar to British Isles, except KGE (2009) [10] and
426 KGE (2012) [11] are single point clusters. Eastern Europe also has clusters similar to British



427 Isles, with the exception that d [2], which is a single point cluster in British Isles, forms a new
428 cluster with M [38] in Eastern Europe.

429 4.2.2 Temperature

430 Compared to precipitation data, temperature data has a lower number of clusters, which can be
431 attributed to the lower variability in temperature data. The clustering of error metrics for British
432 Isles is shown in Fig. 9. For British Isles, 12 clusters are identified, with 5 single point clusters,
433 namely KGE(2009) [10], KGE(2012) [11], MV [22], SGA [35], and R(Spearman) [36]. Similar
434 to precipitation clusters, several error metrics, including ED [7], IRMSE [9], MAE [13], MAPD
435 [15], MASE [16], MSE [23], NRMSE(IQR) [25], NRMSE(mean) [26], NRMSE(range) [27]
436 and RMSE [32] are assigned to the same cluster.

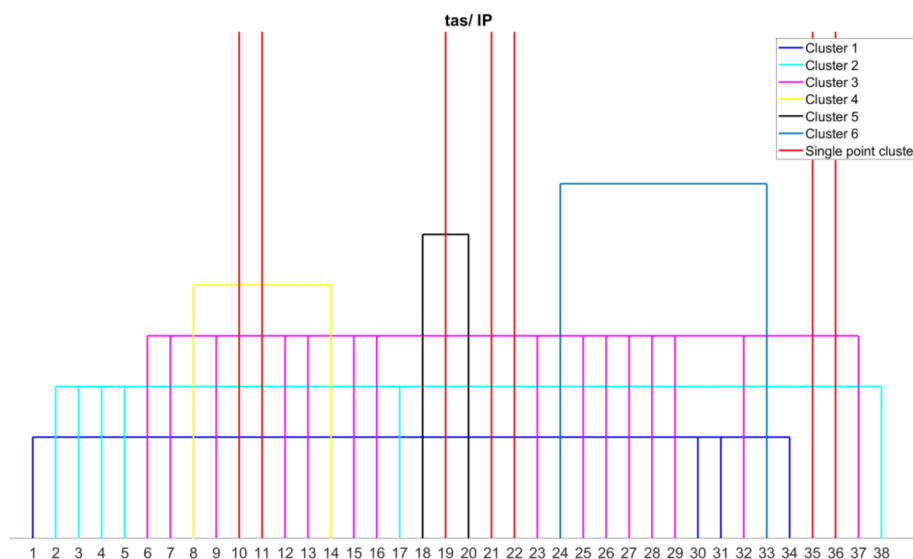


437
438 **Figure 9:** Clustering of error metrics using temperature (tas) data for British Isles (BI) region.
439 Each error metric can be identified by the number using Table S3.

440 The correlation metrics, such as ACC [1], r2 [31], SCO [34], and R(Pearson) [36] belong to
441 the same cluster. France (Fig. S8) and Mid-Europe (Fig. S9) have the same cluster as British
442 Isles for temperature data. For Iberian Peninsula (Fig.10), 13 different clusters are identified,
443 with 7 single point clusters, including MdE [19] and MEE [21] in addition to the 5 single point
444 clusters from British Isles. The remaining clusters are similar to those in British Isles.
445 Mediterranean (Fig. S10) has the same cluster as Iberian Peninsula for temperature data, with
446 13 clusters and 7 single point clusters. Scandinavia (Fig. S11) and Eastern Europe (Fig. S12)



447 have the same number of clusters i.e. 14 clusters. Scandinavia has 8 single point clusters
448 whereas Eastern Europe has 9 single point clusters. Alps (Fig. S13) has 15 clusters with 10
449 single point clusters.



450

451 **Figure 10:** Clustering of error metrics using temperature (tas) data for Iberian Peninsula (IP)
452 region. Each error metric can be identified by the number using Table S3.

453 4.3 Bergen Metrics

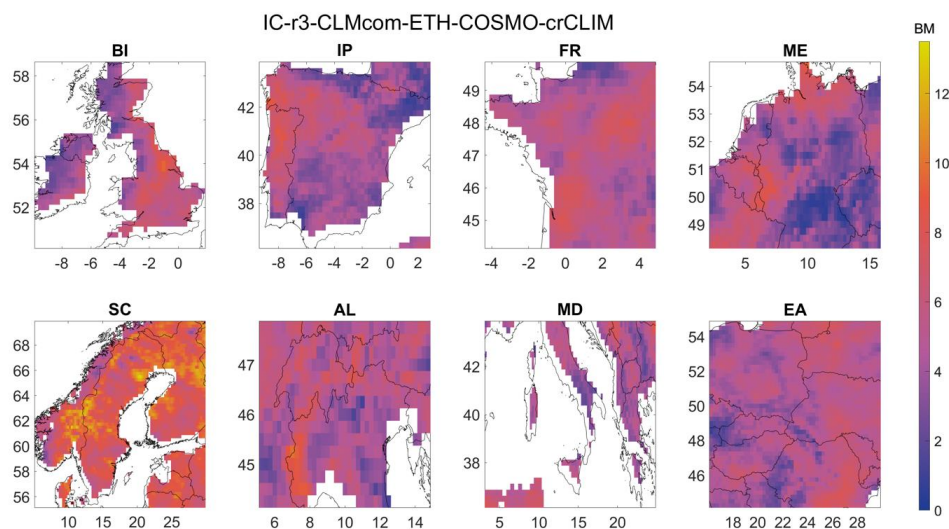
454 A Bergen metric is computed for all eight regions using the respective clusters for both
455 precipitation and temperature. A single metric is chosen from each cluster randomly; Random
456 selection demonstrated no discernible impact on the ranking (see Supplementary Material).
457 Although computed for all 89 regional climate models, this paper focuses on discussing only
458 one climate model for both precipitation and temperature. The CLM Community (CLMCom)
459 regional model from ICHEC-EC-EARTH for r3i1p1 realisation is discussed as it performed
460 best at over 25 grid points in 5 regions and more than 2 grid points in seven regions. For the
461 temperature variable, the CLMCom model form CCCma-CanESM2 model for r1i1p1
462 realisation is discussed, as it performed best at over 25 grid points in seven regions.

463 4.3.1 Precipitation

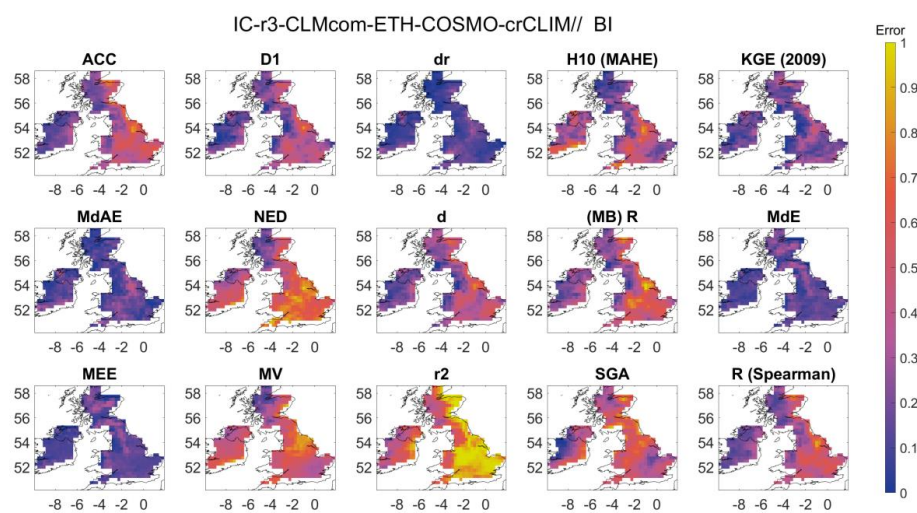
464 A Bergen metric (BM) is used to assess the performance of the CLMCom model for
465 precipitation in all eight different regions. The BM in British Isles region is a composite metric
466 that takes into account 15 different error metrics i.e. ACC, D1, dr, H10(MAHE), KGE(2009),



467 MdAE, NED, d, MB(R), MdE, MEE, MV, r2, SGA, and R(Spearman). Figure 11 provides an
468 overview of the spatial distribution of the BM for all eight regions, while the spatial distribution
469 of each of these metrics is shown in Fig. 12 for the British Isles region.
470 The magnitude of BM ranges from 0 to 13, with a score of 0 indicating good performance by
471 the model. Based on the results, the CLMCom model performed well in the western part of
472 British Isles, as indicated by the BM. This is a result of the good performance of most of the
473 individual metrics that comprise the Bergen Metric. This is shown in Fig. 12. There are some
474 contradictory results from different error metrics in the eastern region. While all 13 metrics
475 indicate good performance, the MV, r2 and NED indicate very bad performance by the model.



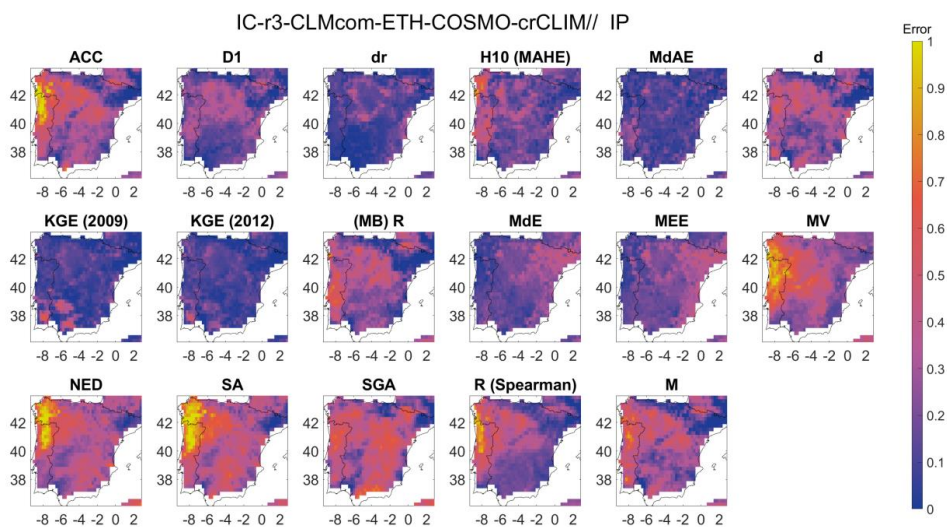
476
477 **Figure 11:** Spatial distribution of Bergen metric using precipitation data for all the eight
478 regions



479

480 **Figure 12:** Spatial distribution of the error metrics used to compute the Bergen metric for
481 precipitation and for British Isles (BI) region. The error metrics have been labelled by the
482 abbreviation and the corresponding error metrics can be identified from Table S3.

483 The use of individual error metrics can provide meaningful insights into the performance of
484 the model in different regions. For example, metrics such as dr, MdAE, MdE, and MEE
485 indicate good performance in the southeastern region, while R(Spearman) indicates bad
486 performance by the CLMCom model which implies that the phase difference is significant
487 between observed and modelled data in this region. It is worth noting that some metrics, such
488 as r2 and R(Spearman), may provide different results even though they share a similar
489 framework. R(Spearman) only tells how well the modelled data follow the observed data while
490 r2 indicate how well the data represents the line of best fit (<https://tinyurl.com/y52r3xed>;
491 <https://tinyurl.com/yk2jmsxt>). Overall, the use of multiple error metrics and the analysis of
492 individual metrics can provide a more comprehensive assessment of the model's performance,
493 particularly in regions where different metrics provide conflicting results.



494

495 **Figure 13:** Spatial distribution of the error metrics used to compute the Bergen metric for
496 precipitation and for Iberian Peninsula (IP) region. The error metrics have been labelled by the
497 abbreviation and the corresponding error metrics can be identified from Table S3.

498 Figure 14 shows a Bergen metric for Iberian Peninsula applied to the CLMCom model, which
499 is based on 17 error metrics obtained from each cluster. These metrics, including ACC, D1, dr,
500 H10 (MAHE), MdAE, d, KGE (2009), KGE (2012), MB (R), MdE, MEE, MV, NED, SA,
501 SGA, R (Spearman) and M, are presented in Fig. 13. The results indicate that the model
502 performs relatively better in the northeast and southeast regions compared to the western region
503 (see Fig. 11), possibly due to the influence of certain metrics such as ACC, R (Spearman), MV,
504 NED, and SA. Additionally, while KGE (2009) and KGE (2012) exhibit similar spatial error
505 patterns, further analysis in the southern region reveals the differences in the magnitude of
506 error. Interestingly, despite their similarity, KGE (2009) and KGE (2012) are classified into
507 different clusters based on a threshold MAE of 5.41, used to determine cluster membership.
508

509 France (Fig. S14), and Mid-Europe (Fig. S15) have the same clusters as the British Isles, and
510 therefore the same error metrics used in British Isles are used to calculate the Bergen metric
511 for France and Mid-Europe. The Bergen metric indicates an average performance of the model
512 for the entire study region of France (see Fig. 11). While r_2 shows a very poor performance of
513 the model for France, MEE metric shows a completely opposite trend, indicating a very good
514 performance of the model. Similar disagreement between r_2 and MEE is also observed in the



515 British Isles. On the other hand, SGA, which compares the shape of the two signals, shows an
516 average performance by the model. In terms of the spatial distribution of error, the Bergen
517 metric shows lower error magnitudes for MEE in the southeast part of the study region.

518 The Bergen metric is also used to assess the performance of the CLMCom model for
519 Scandinavia and Alps using 16 error metrics from each cluster, including ACC, D1, dr, H10
520 (MAHE), MdAE, NED, d, KGE (2009), KGE (2012), MB (R), MdE, MEE, MV, SGA, R
521 (Spearman) and M. The spatial distribution of these metrics is presented in supplementary Fig.
522 S16 (Scandinavia) and Fig. S17 (Alps).

523 Fig. S16 and Fig. 11 suggest that the CLMCom model does not perform well for Scandinavia.
524 However, some error metrics, including dr, MdAE, MdE, and MEE, show good performance
525 in the southern part of the region. Although MdAE, MdE, and MEE are assigned to different
526 clusters, they exhibit similar spatial distributions of error. It is worth noting that despite the
527 similarity, the three error metrics are in different clusters due to their higher MAE between
528 them. For the Alps, the Bergen metric indicates a relatively good performance of the CLMCom
529 model. It can be observed in Fig. S17, all metrics except r2 show good performance for the
530 model.

531 The Mediterranean has the same clusters as the Iberian Peninsula, and the spatial distribution
532 of each metric for the Mediterranean is presented in Fig. S18. The Bergen metric for the
533 CLMCom model suggests an average performance for the entire Mediterranean region. Some
534 of the error metrics, such as KGE (2009), KGE (2012), dr, and MdAE, indicate good model
535 performance. However, metrics such as SGA, SA, and NED, show relatively poor performance
536 of the model.

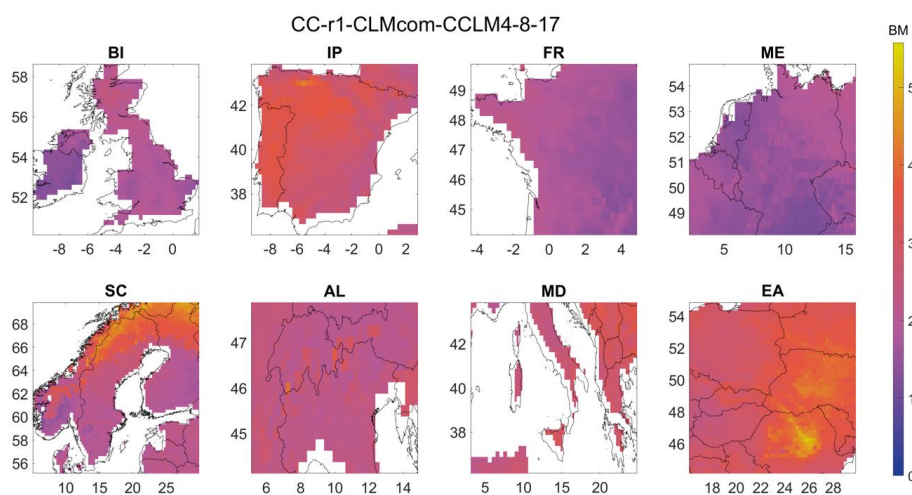
537 For Eastern Europe, the Bergen metric is computed using 14 error metrics from each cluster,
538 as listed: ACC, d, D1, dr, H10(MAHE), KGE(2009), MdAE, NED, MB(R), MdE, MEE, MV,
539 SGA, and R(Spearman). The spatial distribution of each metric is presented in Fig. S19. One
540 notable observation from the figure is the difference between SGA and MEE, which indicates
541 that although the model data has a low bias, the direction of error of the modelled data is
542 completely different from that of the observed data. This insight can be valuable in identifying
543 areas where the model's performance can be improved.

544 **4.3.2 Temperature**

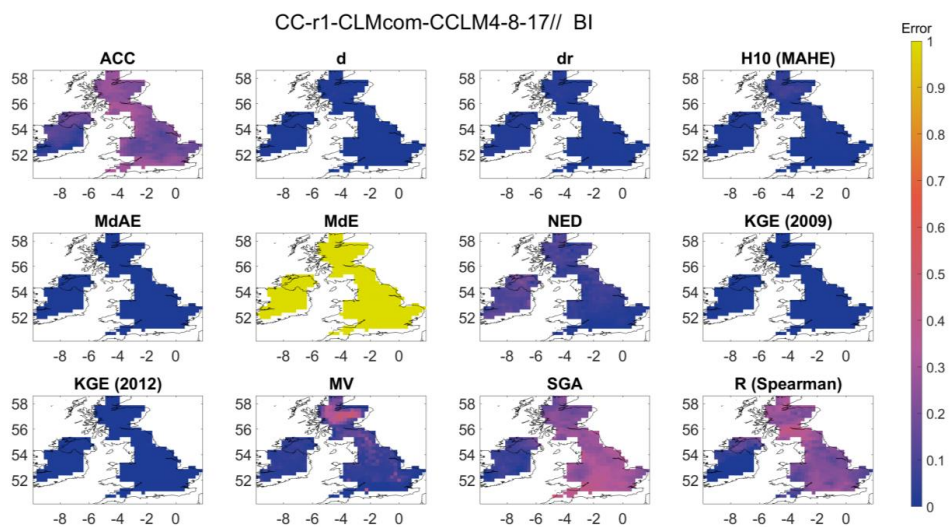
545 For temperature, we focus on the CLM Community (CLMCom) regional model driven by
546 ICHEC-EC-EARTH to demonstrate the application of Bergen metrics for temperature. The
547 spatial distribution of BM is shown in Fig. 14, which indicates average performance by the



548 model, except in certain areas like northern part of Scandinavia, central part of Eastern Europe
549 and western part of Iberian Peninsula, where the performance is bad. The British Isles (Fig.
550 15), France (Fig. S20), and Mid-Europe (Fig. S21) regions have 12 clusters, and 12 error
551 metrics, including ACC, d, dr, H10(MAHE), MdAE, MdE, NED, KGE(2009), KGE(2012),
552 MV, SGA, and R(Spearman) are used to compute the Bergen metric for these regions.

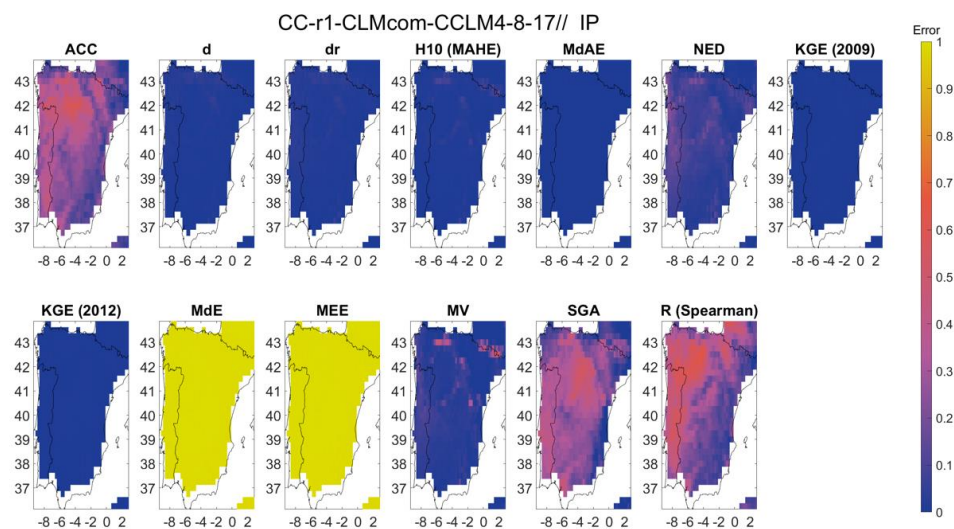


553
554 **Figure 14:** Spatial distribution of Bergen metric using temperature data for all the eight regions
555 The Scandinavia (Fig. S22) and Eastern Europe (Fig. S23) regions have 14 clusters and all the
556 error metrics from British Isles, along with VE and SA, are used to compute the Bergen metric
557 for these regions. The Iberian Peninsula (Fig. 16) and Mediterranean (Fig. S24) regions have
558 the same cluster, with a total of 13 clusters and all the error metrics from British Isles, plus
559 MEE, are used to compute the Bergen metric. The Alps (Fig. S25) region has 15 clusters, with
560 all the error metrics from Scandinavia, including MEE, used to compute the Bergen metric.
561 MdE and MEE consistently indicate very bad model performance for all the regions, while the
562 other metrics indicate relatively good performance. This suggests that the mean and median of
563 the modelled data tend to underestimate/overestimate the observed mean and median,
564 respectively. Histograms in Fig. 17 further investigate this, showing that the error values for
565 ACC are more evenly distributed in the Iberian Peninsula region and close to its ideal point 1,
566 while the source errors for MdE and MEE are concentrated between -0.5 to -1.5, resulting in
567 most of the error values being concentrated between 0.9 to 1 after normalization. The source
568 error represents the distance between the ideal values and actual magnitude after normalization.
569 Similar patterns can be observed in the other regions for temperature .



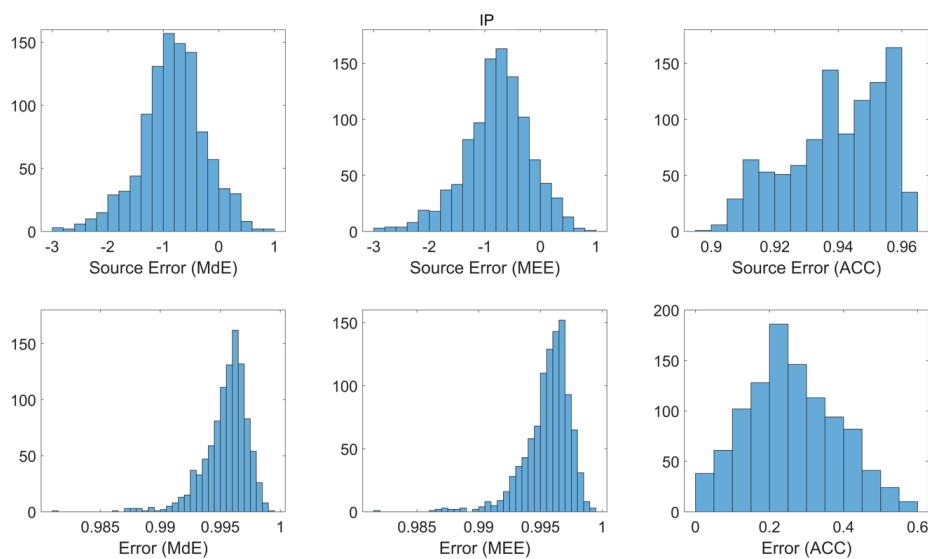
570

571 **Figure 15:** Spatial distribution of the error metrics used to compute the Bergen metric for
 572 temperature and for British Isles (BI) region. The error metrics have been labelled by the
 573 abbreviation and the corresponding error metrics can be identified from Table S3.



574

575 **Figure 16:** Spatial distribution of the error metrics used to compute the Bergen metric for
 576 temperature and for Iberian Peninsula (IP) region. The error metrics have been labelled by the
 577 abbreviation and the corresponding error metrics can be identified from Table S3.



578

579 **Figure 17:** Histogram plot of error and source error for MdE, MEE and ACC for Iberian
580 Peninsula region (IP).

581 5. Conclusions

582 A framework of new error metrics, known as 'Bergen metrics', has been introduced in this study
583 to evaluate the ability of climate models to simulate the observed climate through comparison
584 with a reference field. The proposed metric integrates several error metrics, as described in the
585 results section. To generate a single composite index, the methodology uses a generalized p-
586 norm framework to merge all the error metrics. The research determines that the first norm is
587 the most effective norm to use in the analysis.

588 The study also shows that the number of error metrics used in Bergen Metrics can be reduced
589 using a non-parametric clustering technique. Although several clustering techniques are
590 already available in the literature, they come with certain requirements. Either they require the
591 number of clusters before running the algorithm or information on the class label of the feature
592 vector. The adopted clustering technique tries to identify the natural cluster present in the data.
593 The mean absolute error based on ranking order is used as a dissimilarity index to assign error
594 metrics to different clusters. The technique also has a threshold parameter 5th, 10th and 20th are
595 selected as candidates for threshold parameter and 10th percentile of the D matrix is adopted as
596 a threshold in this study. It is selected because increase in threshold (20th percentile) resulted
597 in increase in MAE and decrease in number of clusters, whereas, decrease in threshold (5th



598 percentile) resulted in decrease in MAE and increase in number of clusters and the study chose
599 a middle ground. However, users can investigate different values of q before choosing the
600 threshold. The clustering technique is compared with the K-means clustering approach and it
601 is found that the non-parametric technique has lower MAE compared to the K-means approach.
602 The clustering is performed for all the eight regions and those are British Isles, Iberian
603 Peninsula, France, Mid-Europe, Scandinavia, Alps, Mediterranean and Eastern Europe. For
604 precipitation, 15, 17, 15, 15, 16, 15, 17, and 14 clusters are obtained for the eight regions,
605 respectively. For temperature, 12, 13, 12, 12, 14, 15, 13, and 14 clusters are obtained for the
606 eight regions, respectively.

607 A single error metric from each cluster can be chosen randomly as a component to be used in
608 the calculation of a Bergen Metric. We have shown that random selection does not have any
609 effect on the ranking order produced by a Bergen Metric. The Bergen Metric which uses the
610 L1 framework is found to be less sensitive to outliers compared to the other norms and more
611 stable in higher dimensional space. Bergen Metrics are a multivariate error functions that can
612 take any number of error metrics of different variables as shown in the last section. It can be
613 further modified for a weighting-based metric that can allow the user to give more weightage
614 to particular metrics depending on the requirement of the study. While some metrics show good
615 performance in certain regions, others indicate poor performance. It is also important to observe
616 how a single metric can influence and change the ranking of climate models. Bergen metrics
617 provide a comprehensive evaluation of the model's performance, which is useful for identifying
618 the strengths and weaknesses of the model in different contexts.

619 Future research should address the sampling uncertainty associated with Bergen metrics. Each
620 data point in time series data has a certain contribution to the total error and if the contribution
621 is not evenly distributed for all the data points, the metric may give biased results. Also, each
622 metric has probabilistic uncertainty associated with it. For example, RMSE works well when
623 the errors are normally distributed and what if the errors are not normally distributed.
624 Discussion on uncertainty may yield useful information that will be helpful in removing the
625 bias from climate models in the future.

626
627
628
629
630
631



632

633

634 **Data and Code availability**

635 The EURO-CORDEX data used in this work are obtained from the Earth System Grid

636 Federation server. The reference precipitation and temperature data is available at

637 <https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-pressure-levels-monthly->

638 [means-preliminary-back-extension?tab=form](https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-pressure-levels-monthly-means-preliminary-back-extension?tab=form)

639 The code for clustering the error metrics is available at [https://github.com/badal01/Error-](https://github.com/badal01/Error-metrics-clustering)

640 [metrics-clustering](https://github.com/badal01/Error-metrics-clustering).

641

642 **Author contributions**

643 AS developed the methodology and performed the formal analysis. PM supervised the research

644 activity planning and execution. AS prepared the first draft of manuscript. All authors

645 contributed to editing and reviewing the manuscript.

646

647 **Competing interests**

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

649

650 **Acknowledgements**

651 The FRONTIER project has received funding from the Research Council of Norway (project

652 number 301777). We thank James Done and Andreas Prein for their advice and critical

653 comments regarding the work.

654

655

656

657

658

659

660

661

662

663

664

665



666 References

- 667 Aggarwal, C. C., Hinneburg, A., & Keim, D. A.: On the surprising behavior of distance metrics in high
668 dimensional space. In *International conference on database theory* (pp. 420-434). Springer, Berlin, Heidelberg,
669 DOI: 10.1007/3-540-44503-X_27, 2001.
- 670 Armstrong, J. S., & Collopy, F.: Error measures for generalizing about forecasting methods: Empirical
671 comparisons. *International journal of forecasting*, 8(1), 69-80, [https://doi.org/10.1016/0169-2070\(92\)90008-W](https://doi.org/10.1016/0169-2070(92)90008-W),
672 1992.
- 673 Baker, N. C., & Taylor, P. C.: A framework for evaluating climate model performance metrics. *Journal of*
674 *Climate*, 29(5), 1773-1782, <https://doi.org/10.1175/JCLI-D-15-0114.1>, 2016
- 675 Bellomo, K., Angeloni, M., Corti, S., & von Hardenberg, J.: Future climate change shaped by inter-model
676 differences in Atlantic meridional overturning circulation response. *Nature Communications*, 12(1), 1-10,
677 <https://doi.org/10.1038/s41467-021-24015-w>, 2021.
- 678 Chai, T., & Draxler, R. R.: Root mean square error (RMSE) or mean absolute error (MAE)?—Arguments against
679 avoiding RMSE in the literature. *Geoscientific model development*, 7(3), 1247-1250,
680 <https://doi.org/10.5194/gmd-7-1247-2014>, 2014.
- 681 Christensen, J. H., & Christensen, O. B.: A summary of the PRUDENCE model projections of changes in
682 European climate by the end of this century. *Climatic change*, 81(Suppl 1), 7-30,
683 <https://doi.org/10.1007/s10584-006-9210-7>, 2007.
- 684 Collier, N., Hoffman, F. M., Lawrence, D. M., Keppel-Aleks, G., Koven, C. D., Riley, W. J., ... & Randerson, J.
685 T.: The International Land Model Benchmarking (ILAMB) system: design, theory, and implementation. *Journal*
686 *of Advances in Modeling Earth Systems*, 10(11), 2731-2754, <https://doi.org/10.1029/2018MS001354>, 2018.
- 687 Flato, G., Marotzke, J., Abiodun, B., Braconnot, P., Chou, S. C., Collins, W., Cox, P., Driouech, F., Emori, S.,
688 Eyring, V., Forest, C., Gleckler, P., Guilyardi, E., Jakob, C., Kattsov, V., Reason, C., & Rummukainen, M.:
689 Evaluation of climate models. In *Climate Change 2013: the physical science basis. Contribution of Working*
690 *Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*, 741-866, [Stocker,
691 T.F., D. Qin, G.-K. Plattner, M. Tignor, S.K. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex and P.M. Midgley
692 (eds.)].Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA, 2013.
- 693 Graham, R. M., Cohen, L., Ritzhaupt, N., Segger, B., Graversen, R. G., Rinke, A., Walden, V.P., Granskog,
694 M.A., & Hudson, S. R.: Evaluation of six atmospheric reanalyses over Arctic sea ice from winter to early
695 summer. *Journal of Climate*, 32(14), 4121-4143, <https://doi.org/10.1175/JCLI-D-18-0643.1>, 2019.
- 696 Gupta, H. V., Kling, H., Yilmaz, K. K., & Martinez, G. F.: Decomposition of the mean squared error and NSE
697 performance criteria: Implications for improving hydrological modelling. *Journal of hydrology*, 377(1-2), 80-91,
698 <https://doi.org/10.1016/j.jhydrol.2009.08.003>, 2009.
- 699 Hartigan, J. A., & Wong, M. A.: Algorithm AS 136: A k-means clustering algorithm. *Journal of the royal*
700 *statistical society. series c (applied statistics)*, 28(1), 100-108, <https://doi.org/10.2307/2346830>, 1979.
- 701 He, X., Lei, X. D., & Dong, L. H.: How large is the difference in large-scale forest biomass estimations based
702 on new climate-modified stand biomass models?. *Ecological Indicators*, 126, 107569,
703 <https://doi.org/10.1016/j.ecolind.2021.107569>, 2021.
- 704 Hu, Z., Chen, X., Zhou, Q., Chen, D., & Li, J.: DISO: A rethink of Taylor diagram. *International Journal of*
705 *Climatology*, 39(5), 2825-2832, <https://doi.org/10.1002/joc.5972>, 2019.
- 706 Hyndman, R. J., & Koehler, A. B.: Another look at measures of forecast accuracy. *International journal of*
707 *forecasting*, 22(4), 679-688, <https://doi.org/10.1016/j.ijforecast.2006.03.001>, 2006.



- 708 IPCC: Climate Change 2021: The Physical Science Basis. Contribution of Working Group I to the Sixth
709 Assessment Report of the Intergovernmental Panel on Climate Change, edited by: Masson-Delmotte, V., Zhai,
710 P., Pirani, A., Connors, S. L., Péan, C., Berger, S., Caud, N., Chen, Y., Goldfarb, L., Gomis, M. I., Huang, M.,
711 Leitzell, K., Lonnoy, E., Matthews, J. B. R., Maycock, T. K., Waterfield, T., Yelekçi, O., Yu, R., & Zhou, B.,
712 Cambridge University Press, https://report.ipcc.ch/ar6/wg1/IPCC_AR6_WGI_FullReport.pdf, 2021a.
- 713 IPCC: Summary for Policymakers, in: Climate Change 2021: The Physical Science Basis. Contribution of
714 Working Group I to the Sixth Assessment Report of the Intergovernmental Panel on Climate Change, edited by:
715 Masson-Delmotte, V., Zhai, P., Pirani, A., Connors, S. L., Péan, C., Berger, S., Caud, N., Chen, Y., Goldfarb,
716 L., Gomis, M. I., Huang, M., Leitzell, K., Lonnoy, E., Matthews, J. B. R., Maycock, T. K., Waterfield, T.,
717 Yelekçi, O., Yu, R., & Zhou, B., Cambridge University Press,
718 https://www.ipcc.ch/report/ar6/wg1/downloads/report/IPCC_AR6_WGI_SPM.pdf, 2021b.
- 719 Jackson, E. K., Roberts, W., Nelsen, B., Williams, G. P., Nelson, E. J., & Ames, D. P.: Introductory overview:
720 Error metrics for hydrologic modelling—A review of common practices and an open source library to facilitate
721 use and adoption. *Environmental Modelling & Software*, 119, 32-48,
722 <https://doi.org/10.1016/j.envsoft.2019.05.001>, 2019.
- 723 Kalmár, T., Pieczka, I., & Pongrácz, R.: A sensitivity analysis of the different setups of the RegCM4.5 model
724 for the Carpathian region. *International Journal of Climatology*, 41, E1180-E1201,
725 <https://doi.org/10.1002/joc.6761>, 2021.
- 726 Kling, H., Fuchs, M., & Paulin, M.: Runoff conditions in the upper Danube basin under an ensemble of climate
727 change scenarios. *Journal of hydrology*, 424, 264-277, <https://doi.org/10.1016/j.jhydrol.2012.01.011>, 2012.
- 728 Kotlarski, S., Keuler, K., Christensen, O. B., Colette, A., Déqué, M., Gobiet, A., Goergen, K., Jacob, D., Lüthi,
729 D., van Meijgaard, E., Nikulin, G., Schär, C., Teichmann, C., Vautard, R., Warrach-Sagi, K., & Wulfmeyer, V.:
730 Regional climate modeling on European scales: a joint standard evaluation of the EURO-CORDEX RCM
731 ensemble. *Geoscientific Model Development*, 7(4), 1297-1333. <https://doi.org/10.5194/gmd-7-1297-2014>, 2014.
- 732 Liemohn, M. W., Shane, A. D., Azari, A. R., Petersen, A. K., Swiger, B. M., & Mukhopadhyay, A.: RMSE is
733 not enough: Guidelines to robust data-model comparisons for magnetospheric physics. *Journal of Atmospheric
734 and Solar-Terrestrial Physics*, 218, 105624, <https://doi.org/10.1016/j.jastp.2021.105624>, 2021.
- 735 Mirkes, E. M., Allohibi, J., & Gorban, A.: Fractional norms and quasinorms do not help to overcome the curse
736 of dimensionality. *Entropy*, 22(10), 1105, <https://doi.org/10.48550/arXiv.2004.14230>, 2020.
- 737 Mooney, P. A., Rechid, D., Davin, E. L., Katragkou, E., de Noblet-Ducoudré, N., Breil, M., Cardoso, R. M.,
738 Daloz, A. S., Hoffmann, P., Lima, D. C. A., Meier, R., Soares, P. M. M., Sofiadis, G., Strada, S., Strandberg, G.,
739 Toelle, M. H., & Lund, M. T.: Land-atmosphere interactions in sub-polar and alpine climates in the CORDEX
740 Flagship Pilot Study Land Use and Climate Across Scales (LUCAS) models – Part 2: The role of changing
741 vegetation, *The Cryosphere*, 16, 1383–1397, <https://doi.org/10.5194/tc-16-1383-2022>, 2022
- 742 Murphy, A. H.: Skill scores based on the mean square error and their relationships to the correlation
743 coefficient. *Monthly weather review*, 116(12), 2417-2424. [https://doi.org/10.1175/1520-0493\(1988\)116%3C2417:SSBOTM%3E2.0.CO;2](https://doi.org/10.1175/1520-0493(1988)116%3C2417:SSBOTM%3E2.0.CO;2), 1988
- 745 Nash, J. E., & Sutcliffe, J. V.: River flow forecasting through conceptual models part I—A discussion of
746 principles. *Journal of hydrology*, 10(3), 282-290. [http://dx.doi.org/10.1016/0022-1694\(70\)90255-6](http://dx.doi.org/10.1016/0022-1694(70)90255-6), 1970.
- 747 Pachepsky, Y. A., Martinez, G., Pan, F., Wagener, T., & Nicholson, T.: Evaluating hydrological model
748 performance using information theory-based metrics. *Hydrology and Earth System Sciences Discussions*, 1-24,
749 <https://doi.org/10.5194/hess-2016-46>, 2016.
- 750 Pitman, J.: Exchangeable and partially exchangeable random partitions. *Probability theory and related
751 fields*, 102(2), 145-158, <https://doi.org/10.1007/BF01213386>, 1995.



- 752 Randall, D. A., Wood, R. A., Bony, S., Colman, R., Fichet, T., Fyfe, J., Kattsov, V., Pitman, A., Shukla, J.,
753 Srinivasan, J., Stouffer, R.J., Sumi, A. & Taylor, K. E.: Climate models and their evaluation. In *Climate*
754 *Change 2007: The physical science basis. Contribution of Working Group I to the Fourth Assessment Report of*
755 *the IPCC (FAR)* (pp. 589-662). Cambridge University Press, 2007.
- 756 Reich, N. G., Lauer, S. A., Sakrejda, K., Iamsrithaworn, S., Hinjoy, S., Suangtho, P., Suthachana, S., Clapham,
757 H.E., Salje, H., Cummings, D.A. & Lessler, J.: Challenges in real-time prediction of infectious disease: a case
758 study of dengue in Thailand. *PLoS neglected tropical diseases*, 10(6), e0004761,
759 <https://doi.org/10.1371/journal.pntd.0010883>, 2016.
- 760 Richter, J. H., Butchart, N., Kawatani, Y., Bushell, A. C., Holt, L., Serva, F., Anstey, J., Simpson, I.R., Osprey,
761 S., Hamilton, K., Braesicke, P., Cagnazzo, C., Chen C. C., Garcia, R. R., Gray, L. J., Kerzenmacher, T., Lott, F.,
762 McLandress, C., Naoe, H., Scinocca, J., Stockdale, T. N., Versick, S., Watanabe, S., Yoshida, K., & Yukimoto,
763 S.: Response of the quasi-biennial oscillation to a warming climate in global climate models. *Quarterly Journal*
764 *of the Royal Meteorological Society*, 148(744), 1490-1518, <https://doi.org/10.1002/qj.3749>, 2022.
- 765 Rupp, D. E., Abatzoglou, J. T., Hegewisch, K. C., & Mote, P. W.: Evaluation of CMIP5 20th century climate
766 simulations for the Pacific Northwest USA. *Journal of Geophysical Research: Atmospheres*, 118(19), 10-884,
767 <https://doi.org/10.1002/jgrd.50843>, 2013.
- 768 Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram. *Journal of Geophysical*
769 *Research: Atmospheres*, 106(D7), 7183-7192, <https://doi.org/10.1029/2000JD900719>, 2001.
- 770 van Noije, T., Bergman, T., Le Sager, P., O'Donnell, D., Makkonen, R., Gonçalves-Ageitos, M., M., Döschner,
771 R., Fladrich, U., von Hardenberg, J., Keskinen, J.-P., Korhonen, H., Laakso, A., Myriokefalitakis, S., Ollinaho,
772 P., Pérez García-Pando, C., Reerink, T., Schrödner, R., Wyser, K., & Yang, S.: EC-Earth3-AerChem: a global
773 climate model with interactive aerosols and atmospheric chemistry participating in CMIP6. *Geoscientific Model*
774 *Development*, 14(9), 5637-5668, <https://doi.org/10.5194/gmd-14-5637-2021>, 2021.
- 775 Watt-Meyer, O., Brenowitz, N. D., Clark, S. K., Henn, B., Kwa, A., McGibbon, J., Perkins, W. A. &
776 Bretherton, C. S.: Correcting weather and climate models by machine learning nudged historical
777 simulations. *Geophysical Research Letters*, 48(15), e2021GL092555, <https://doi.org/10.1029/2021GL092555>,
778 2021.
- 779 Weber, R., Schek, H. J., & Blott, S.: A quantitative analysis and performance study for similarity-search
780 methods in high-dimensional spaces. In *VLDB*, 98, 194-205, 1998.
- 781 Węglarczyk, S.: The interdependence and applicability of some statistical quality measures for hydrological
782 models. *Journal of Hydrology*, 206(1-2), 98-103. [https://doi.org/10.1016/S0022-1694\(98\)00094-8](https://doi.org/10.1016/S0022-1694(98)00094-8), 1998.
- 783 Wehner, M., Lee, J., Risser, M., Ullrich, P., Gleckler, P., & Collins, W. D.: Evaluation of extreme sub-daily
784 precipitation in high-resolution global climate model simulations. *Philosophical Transactions of the Royal*
785 *Society A*, 379(2195), 20190545, <https://doi.org/10.1098/rsta.2019.0545>, 2021.
- 786 Wilcox, R. H.: Adaptive control processes—A guided tour, by Richard Bellman, Princeton University Press,
787 Princeton, New Jersey, 1961, 255 pp., \$6.50. *Naval Research Logistics Quarterly*, 8(3), 315-316,
788 <https://www.jstor.org/stable/j.ctt183ph6v>, 1961.
- 789 Willmott, C. J., & Matsuura, K.: Advantages of the mean absolute error (MAE) over the root mean square error
790 (RMSE) in assessing average model performance. *Climate research*, 30(1), 79-82,
791 <https://www.jstor.org/stable/24869236>, 2005.
- 792 Wood, R. R., Lehner, F., Pendergrass, A. G., & Schlunegger, S.: Changes in precipitation variability across time
793 scales in multiple global climate model large ensembles. *Environmental Research Letters*, 16(8), 084022.
794 <https://doi.org/10.1088/1748-9326/ac10dd>, 2021.



795 Yang, J., Ren, J., Sun, D., Xiao, X., Xia, J. C., Jin, C., & Li, X.: Understanding land surface temperature impact
796 factors based on local climate zones. *Sustainable Cities and Society*, 69, 102818,
797 <https://doi.org/10.1016/j.scs.2021.102818>, 2021.

798