

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

3 Alok K. Samantaray<sup>1,2</sup>, Priscilla A. Mooney<sup>1,2</sup>, Carla A. Vivacqua<sup>3</sup>

4 <sup>1</sup>Norwegian Research Centre (Norce), Norway

5 <sup>2</sup>Bjerknes Centre for Climate Research, Norway

6 <sup>3</sup>Federal University at Rio Grande do Norte, Brazil

7  
8 Corresponding author:

9 Alok Kumar Samantaray

10 Jahnebakken 5, 5007 Bergen, Norway

11 Email: [asam@norceresearch.no](mailto:asam@norceresearch.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 **Abstract**

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

72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83

## 84        **1. Introduction**

85    Climate models are important tools for predicting and understanding climate change, and  
86    climate processes (Kotlarski et al., 2014; IPCC, 2021a; IPCC, 2021b; Mooney et al., 2022). In  
87    the context of climate studies, climate model evaluation is essential for identifying models that  
88    poorly simulate the climate system, and for ranking of climate models (Randall et al., 2007;  
89    Flato et al., 2013). The main purpose of climate model evaluation is twofold; firstly, to ensure  
90    that the models are reproducing key aspects of the climate system and secondly to understand  
91    the limitations of climate projections from the models. This ensures proper interpretation and  
92    application of climate models and any climate projections produced by them. The performance  
93    of climate models is quantified by different error metrics such as root mean square error, and  
94    bias, which assess the agreement between the climate model data and reference data (e.g.,  
95    gridded observational products, station data, reanalyses, or satellite observations). As the  
96    number of climate models has increased, the study of error metrics has become increasingly  
97    important. There are several error metrics available to evaluate the performance of climate  
98    models (Jackson et al., 2019), and the selection of an appropriate metric remains a topic of  
99    debate in the literature. For instance, Willmott & Matsuura (2005) advocate for mean absolute  
100    error (MAE) over root mean squared error (RMSE), as the latter is not an effective indicator of  
101    average model performance. In contrast, Chai & Draxler (2014) contend that RMSE is superior  
102    to MAE when errors follow a Gaussian distribution.

103    Different error metrics are available in the literature, and each has a specific framework  
104    according to its purpose (Rupp et al., 2013; Pachepsky et al., 2016; Baker & Taylor, 2016;  
105    Collier et al., 2018; Jackson et al., 2019). For example, root mean square error compares the  
106    amplitude difference between modelled and reference data, while the correlation coefficient  
107    compares the phase difference between modelled and reference data. Depending on the specific  
108    error, the error metrics can be categorised into different classes; the most popular classes are  
109    accuracy, precision, and association. Accuracy measures the degree of similarity between  
110    climate model data and reference data. An extremely high accuracy indicates that the model  
111    has less error magnitude of any type and testing the model with other error metrics adds little  
112    value (Liemohn et al., 2021). However, if a model has moderate to low accuracy, testing the  
113    model with other metrics can reveal other similarities and dissimilarities between model data  
114    and reference data. Root mean square error and mean square error are the most used accuracy  
115    metrics to evaluate climate models (Watt-Meyer et al., 2021; Wehner et al., 2021; He et al.,  
116    2021), even though the metrics cannot reveal whether the model is under or over-predicting  
117    the observations. Precision metrics quantify the degree of similarity in the spread of the data.

118 A robust and commonly used metric for assessing the precision of model data is the ratio or  
119 difference of standard deviation between modelled data and reference data (van Noije et al.,  
120 2021; Wood et al., 2021; Wehner et al., 2021). Finally, association metrics measure the degree  
121 of the phase difference between modelled data and observed data. Phase difference is important  
122 in climate studies as it affects the initiation and termination time of a season of climate  
123 variables. One metric that is extensively used to measure the association is the correlation  
124 coefficient (Richter et al., 2022; Bellomo et al., 2021; Yang et al., 2021). Liemohn et al. (2021)  
125 has described various other major categories of metrics and they suggest that assessment of  
126 models should not be restricted to one or two error metrics. Interested readers can follow the  
127 citations to read in detail about the discussed metrics.

128 In addition to this, researchers have employed various characteristics of climatic parameters as  
129 measures to assess and compare climate models with observed datasets. Metrics encompassing  
130 the frequency of days with precipitation over 1 mm and over 15 mm, the 90% quantile of the  
131 frequency distribution, and the maximum number of consecutive dry days, along with  
132 parameters such as daily mean, daily maximum, daily minimum, yearly maximum, length of  
133 the frost-free period, growing degree days ( $> 5^{\circ}\text{C}$ ), cooling degree days ( $> 22^{\circ}\text{C}$ ), heating  
134 degree days ( $< 15.5^{\circ}\text{C}$ ), days with RR ( $> 99$ th percentile of daily amounts for all days), ratio  
135 of spatial variability, pattern correlation, ratio of interannual variability, temporal correlation  
136 of interannual variability, number of summer days, number of frost days, consecutive dry days,  
137 and ratio of yearly amplitudes, have been utilized for the validation of Euro-CORDEX data  
138 (Kotlarski et al., 2014; Giot et al., 2016; Smiatek et al., 2016; Torma, 2019; Vautard et al.,  
139 2021). Other studies have employed the empirical orthogonal functions (Rasmus et al., 2023),  
140 structural similarity index metric (Wang & Bovik, 2002), fractions skill score (Roberts & Lean,  
141 2008), spatial pattern efficiency metric (Dembélé et al., 2020), spatial efficiency metric  
142 (Demirel, 2018; Ahmed et al., 2019) and probability distribution function (Perkins et al., 2007;  
143 Boberg et al., 2009; Boberg et al., 2010; Masanganise et al., 2014) to evaluate climate models.  
144 There are several composite error metrics that use the modified framework of other metrics to  
145 compute the error magnitude. A widely used example of this is the Taylor diagram (Taylor,  
146 2001), which incorporates correlation, root mean square deviation and ratio of standard  
147 deviation. A distinguishing feature of the Taylor Diagram is its ability to graphically evaluate  
148 the model performance. Another popular example is the Nash-Sutcliffe Efficiency (NSE; Nash  
149 & Sutcliffe, 1970) which is a normalised form of the mean squared error to evaluate and predict  
150 the model streamflow data. Later, it was observed that NSE can be decomposed into three

151 components which are the functions of correlation, bias and standard deviation (Murphy, 1988;  
152 Weglarczyk, 1998). Other similar scores include the Kling-Gupta (K-G) efficiency (Gupta et  
153 al., 2009) which is a function of three components: ratio of model mean to observed mean, the  
154 ratio of model standard deviation to observed standard deviation and correlation coefficient.  
155 The study of Gupta et al. (2009) argued the NSE, which has a bias component normalised by  
156 the standard deviation of the reference data, will have a low weight on the bias component if  
157 the reference data has high variability. The modified Kling-Gupta efficiency developed by  
158 Kling et al. (2012) involves the ratio of covariance instead of the ratio of standard deviation.  
159 Both K-G efficiency and modified K-G efficiency use Euclidean distance as a basis to calculate  
160 the error magnitude of the model and the study argued that instead of finding a corrected NSE  
161 criterion, the whole problem can be viewed from the multi-objective perspective where the  
162 three error components can be used as separate criteria to be optimised. It identifies the best  
163 models by calculating the Euclidean distance from the ideal point and then finding the model  
164 with the shortest distance. The ideal value of an error metric is obtained when the model exactly  
165 simulates the observed data. The Euclidean distance is also used by Hu et al. (2019) to develop  
166 the DISO metric that incorporates correlation coefficient, absolute error and root mean squared  
167 error. The study of Hu et al. (2019) also argues that accuracy (root mean square error), bias  
168 (absolute error) and association (correlation coefficient) are the three major error classes based  
169 on which a model should be assessed and evaluating a model using a single error metric may  
170 lead to ill-informed results. The study pointed out a few limitations of the Taylor diagram such  
171 as quantification of error magnitude and low sensitivity to small error differences by the  
172 diagram. In a comparative study, Kalmár et al. (2021) found no substantial difference between  
173 DISO index and the Taylor diagram. However, based on quantification of error magnitude,  
174 DISO index can be helpful.

175 The Euclidean distance framework has found increasing use in various fields, serving as an  
176 error function or metric in applications like model evaluation, parameter optimization, and  
177 classification problems. In essence, it calculates the straight-line distance between two points  
178 in the space, known as Euclidean distance. The Euclidean distance is essentially the second  
179 norm of a vector. Equation 1 represents the generalized form of the p-norm in an n-dimensional  
180 vector space, where  $x_i$  is the vector. When p is set to 2, it transforms into the Euclidean norm.  
181 In the context of time series data, if the vector  $(x_i)$  represents the difference between observed  
182 data  $(u_i)$  and model data  $(v_i)$  i.e.,  $x_i = u_i - v_i$ , then d is termed the Euclidean distance metric.

183 Here,  $i$  represents the time series data. It's important to note that root mean squared error and  
 184 mean squared error are different variants of the Euclidean distance metric.  
 185 Furthermore, if the vector represents the difference between error metrics (correlation  
 186 coefficient [ $u_1$ ], absolute error [ $u_2$ ] and root mean squared error [ $u_3$ ]) and their ideal values  
 187 ( $v_{1,3}$ ), then  $d$  is referred to as the DISO index. In summary, the Euclidean distance framework  
 188 offers a versatile approach applicable to various scenarios, providing valuable insights through  
 189 different metrics and indices. A disadvantage of the Euclidean distance is that it suffers the  
 190 curse of dimensionality (Mirkes et al., 2020; Weber et al., 1998) i.e. Euclidean distance as a  
 191 dissimilarity index becomes less efficient as dimension increases. In this study, we assess the  
 192 effect of the norm order on the overall error. We use different measures such as the contribution  
 193 of outliers to the overall error, the difference between the maximum and minimum distances,  
 194 and the average distances to compare different norms.

$$195 \quad d_n(u, v) = (\sum_{i=1}^n |x_i(u_i, v_i)|^p)^{1/p} \quad (1)$$

196 This study has the following objectives:

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

## 202 **2. Data and Study area**

203 We focus on Europe due to the widespread availability of a large ensemble of high resolution  
 204 ( $0.11^\circ$ ) regional climate simulations. In this study, we use 89 regional climate model (RCM)  
 205 simulations from Euro-CORDEX to study the behaviour of different error metrics. The Euro-  
 206 CORDEX dataset provides both precipitation and temperature data at  $0.11^\circ$  grid resolution.  
 207 The monthly data from 1975 to 2005, which is available in all the RCM simulations, have been  
 208 used to calculate the index. Supplementary Table S1 provides an overview of the global climate  
 209 models (GCMs) downscaled by the different RCMs. Supplementary Table S2 provides an  
 210 overview of the RCMs and assigns a number (Column 1) to each RCM which is used to identify  
 211 RCMs in plots that have limited space for labels.

212 For reference data, both precipitation and temperature data are obtained from the E-OBS  
 213 dataset. The study utilized the  $0.25^\circ$  grid resolution dataset to meet the specific requirements  
 214 of the project. However, users can choose datasets of different resolutions based on their study

215 needs for climate model validation. To facilitate the comparison of model data with the  
 216 reference data, all datasets need to be on a common grid. In this study, we remapped the RCM  
 217 data onto the coarser 0.25° grid of E-OBS.

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

### 221 3. Methodology

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

#### 224 3.1 Error metrics

225 Error metrics play a crucial role in climate change studies, serving as essential tools to quantify  
 226 the disparities between modelled and reference data over time series. Each error metric is  
 227 designed to capture specific aspects of the relationship between model data and reference data,  
 228 as discussed in the introduction section. To gain insight into the performance of error metrics,  
 229 we have analysed Euro-CORDEX precipitation data and examined the differences in ranking  
 230 of 89 GCM-driven regional climate simulations using 38 error metrics. The list of error metrics  
 231 is provided in Table S3 and the details of all 38 error metrics have been provided in Jackson et  
 232 al., (2019). All 89 models are ranked based on their performance using the 38 error metrics.  
 233 The average ( $r_{M,mean}$ ; Equation 2) and maximum ( $r_{M,max}$ ; Equation 3) rank differences are  
 234 then calculated at each grid point. The former is the mean of all the pairwise rank differences,  
 235 while the latter is the maximum of all the pairwise rank differences. These calculations allow  
 236 us to understand the performance of different error metrics and the extent of the disparity in  
 237 ranking of the climate models.

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

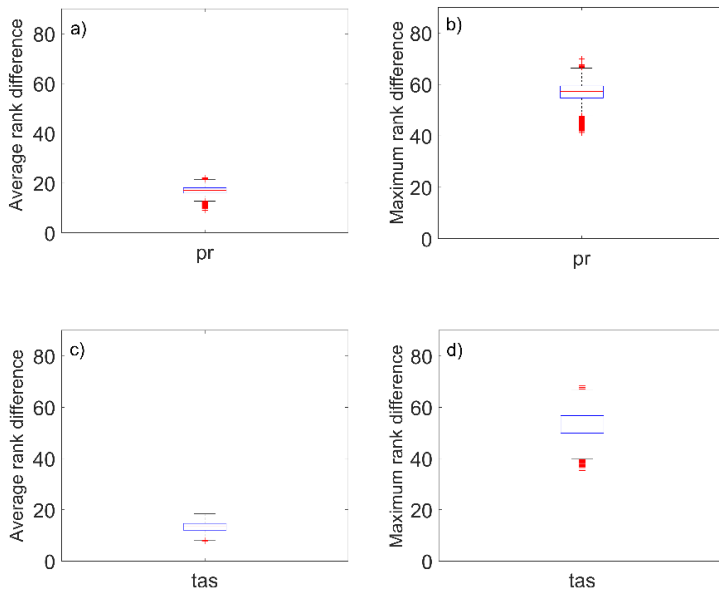
Number	Climate model	Ranking order (RO) by $i$ th error metric ( $E_i$ )	Ranking order (RO) by $k$ th error metric ( $E_k$ )
1	M1	3	2
2	M2	1	3
3	M3	2	1

239

240 
$$r_{M,mean} = \mu_g(R_{M,k} - R_{M,i}) \quad (2)$$

241 
$$r_{M,max} = \max_g (R_{M,k} - R_{M,i}) \quad (3)$$

242  $R_{M,k}$  and  $R_{M,i}$  are the rank assigned to model M by the  $k$ th and  $i$ th error metric, respectively.  
 243 We have provided Table 1 as an example for better understanding of the notations. If there are  
 244 three climate models (M1, M2 and M3) as shown in Table 1, all the models have been assigned  
 245 to a number (first column) and the order must not change throughout the study.  $R_{M,k}$  and  $R_{M,i}$   
 246 for model M1 are 2 and 3, respectively.  $k$  varies from 1 to  $N_E-1$  and  $i$  varies from  $k+1$  to  $N_E$ ,  
 247 where  $N_E$  is the total number of error metrics. The difference in ranking is calculated for all  
 248 possible combinations of error metrics.  $\mu_g()$  and  $\max_g()$  are the mean and maximum operator,  
 249 respectively, which is applied across all the grid points ( $g:1,2,\dots,g_d$ ).  $g_d$  is the total number of  
 250 grid points which is 11370 in this study. Figure 1 demonstrates that different error metrics used  
 251 to assess climate models result in significantly different ranking orders. The average of  $r_{M,mean}$   
 252 across all the grid point varies from 16 to 26 whereas the average of  $r_{M,max}$  varies from 40 to  
 253 70. The results indicate significant differences in the ranking of the climate models by different  
 254 error metrics. The disparity in ranking order may be due to the distinctive error targeted by  
 255 each metrics as discussed in the introduction section.



256 **Figure 1:** Box plot of average rank difference (first column [a, c]) and maximum rank  
 257 difference (second column; [b, d]) for precipitation (Pr; first row [a, b]) and temperature (T;  
 258 second row [c, d]) over all the grid points in European region  
 259

260 This study assumes that all the errors are important and that it may be necessary to evaluate  
 261 model performance using multiple metrics. To achieve independence among the metrics, the  
 262 study has attempted to cluster the error metrics based on model performance. This classification  
 263 would enable different clusters to have unique characteristics, and metrics within the same



264 cluster would produce similar results, whereas those from different clusters would yield  
265 different ranking orders. In summary, the study proposes that using multiple error metrics and  
266 clustering them based on performance could improve the understanding and  
267 comprehensiveness of climate model analysis.

### 268 **3.2 Clustering of error metrics**

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

278 Similar to the rank difference explained in the previous section, the MAE ( $RO_i, RO_k$ ) between  
279 the ranking order produced by two error metrics is computed. RO is the ranking order and it  
280 can be calculated by assigning the climate models to a number. For example, the ranking order  
281 ( $RO_i$ ) by  $i$ th error metric and the ranking order ( $RO_k$ ) by  $k$ th error metric are [3, 1, 2] and [2,  
282 3, 1], respectively in Table 1. The MAE values are calculated for all possible combinations of  
283 error metrics in a particular cluster and the maximum of the MAE values is used to compare it  
284 to the threshold value. The exercise is repeated for all the clusters ( $N_C$ ) available at that time.  
285 The number of clusters ( $N_C$ ) and the number of error metrics in each cluster ( $N_{CE}$ ) are updated  
286 for each iteration (i) and if the criteria is not satisfied, then a new cluster is formed using that  
287 error metric. The whole exercise is repeated till all the error metrics ( $N_E$ ) gets assigned to a  
288 cluster.

289 The threshold value is defined as  $q$ th percentile of a column matrix D where D is the collection  
290 of MAE values for all possible combinations of error metrics at all the grid points in a region.  
291 In this study,  $q$  has been assigned the value of 10 and the sensitivity of  $q$  is discussed in the  
292 results section.

293

```

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

```

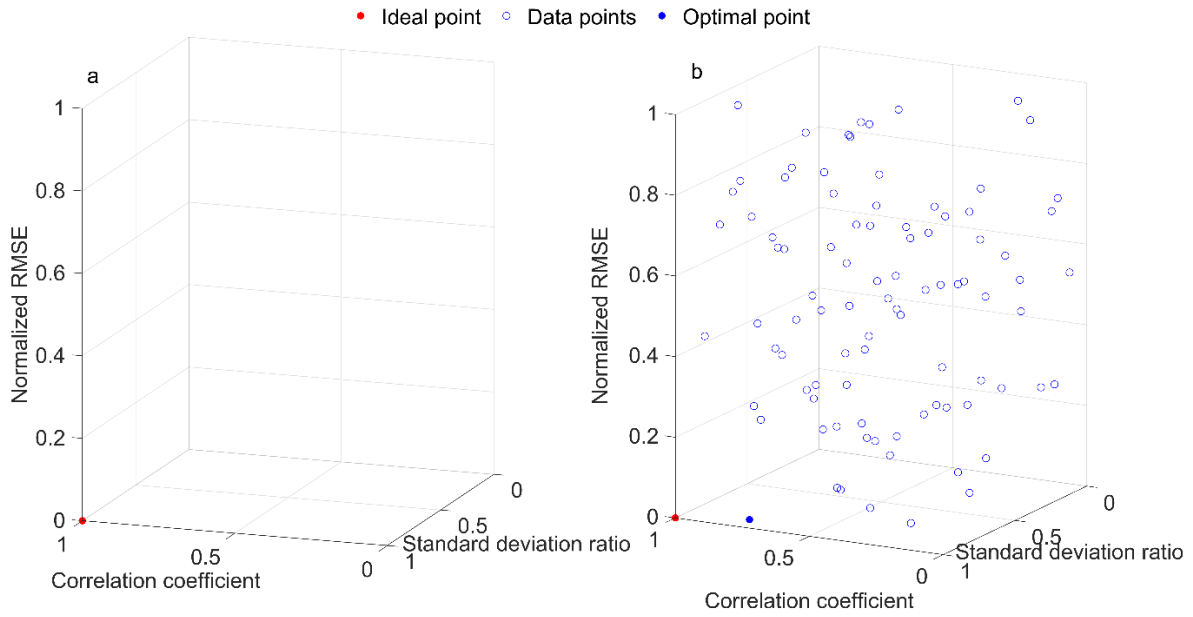
294

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

### 296 3.3 Proposed metric- The Bergen Metrics

297 The clustering of error metrics guarantees that metrics in different groups produce distinct  
 298 ranking orders, implying that each group targets different errors. One of the objectives of this  
 299 study is to integrate different errors and create a composite error to obtain a single value. One  
 300 potential solution is to use the Euclidean distance approach with different error metrics as  
 301 different dimensions in the Euclidean space. To illustrate this, we employed three widely used  
 302 error metrics: Normalized Root Mean Square Error (RMSE), Standard Deviation ratio (SD)  
 303 and correlation coefficient. In the Euclidean space, an ideal model that predicts the climate  
 304 variable as accurately as the observed data would have values of 1, 1, and 0 for correlation  
 305 coefficient, Standard Deviation ratio, and normalized RMSE, respectively. The coordinates of  
 306 an ideal model in the Euclidean space would be (1, 1, 0), as represented by the red point in Fig.  
 307 3a. Since different models have unique coordinates based on the three metrics, these  
 308 coordinates serve as possible solutions to determine the best model. If a decision is required,  
 309 one approach could be to calculate the Euclidean distance from the ideal point to all points and  
 310 select the point with the shortest distance (Equation 4). The model that is closest to the ideal  
 311 point, indicated by the optimal point in Fig.3b, can be considered as the best model.

$$312 \quad ED \text{ Metric} = \sqrt{(1 - \text{Correlation coefficient})^2 + (1 - \text{Standard deviation ratio})^2 + (0 - RMSE)^2} \quad (4)$$

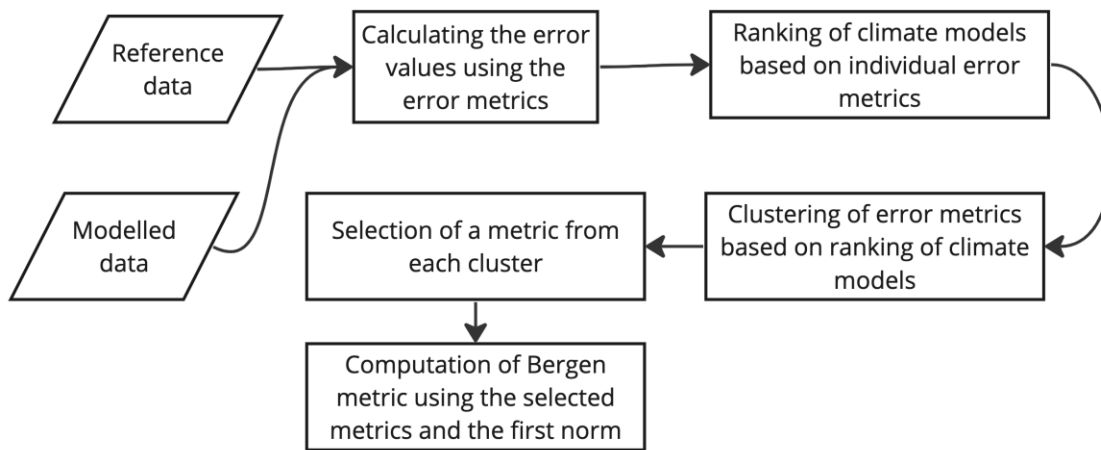


313

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

316

317 The Euclidian distance has several benefits that make it a popular metric, primarily its  
 318 simplistic framework. However, it also has some drawbacks. The Euclidian distance, also  
 319 known as L2 norm, is less effective in higher dimensional spaces, which can lead to instability  
 320 when additional error metrics are added (Weber et al., 1998; Aggarwal et al., 2001). To mitigate  
 321 this issue, recent research has focused on the use of L1 norms, such as relative mean absolute  
 322 error and mean absolute scaled error, which have become more popular than L2 norms like  
 323 mean squared error. This approach reduces the impact of outliers in the data (Armstrong &  
 324 Collopy, 1992; Hyndman and Koehler, 2006). Reich et al. (2016) found that relative MAE,  
 325 based on an L1 norm, is advantageous in assessing prediction models. This study proposes the  
 326 a new metrics called the Bergen Metrics (BM) which is a generalised p-norm framework to  
 327 evaluate climate models.



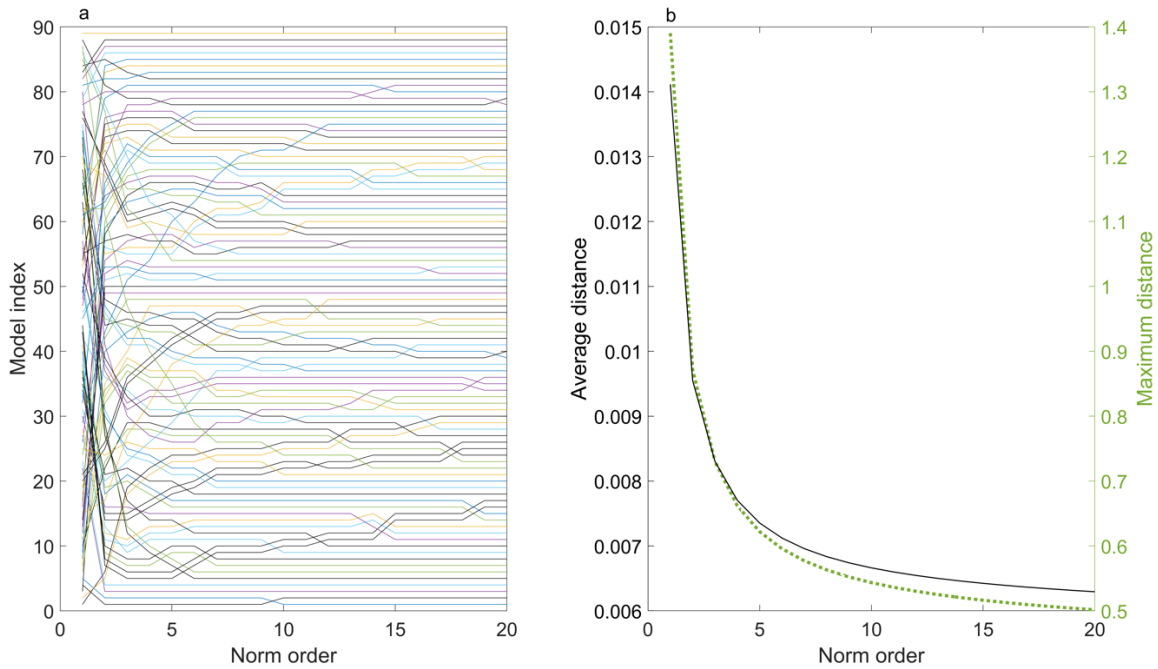
328

329 **Figure 4:** The flowchart for the calculation of Bergen metric

330 A case study has been conducted to understand the impact of different p norms on the ranking  
 331 order of climate models. For this, five error metrics - RMSE, bias, correlation coefficient,  
 332 standard deviation ratio, and mean ratio - have been considered (Equation 5) and the error  
 333 metrics are normalised using model data. A flowchart has been provided to illustrate the various  
 334 steps involved in calculating the Bergen metric (Fig. 4). It is important to note that equation 5  
 335 serves as an illustration of Bergen metrics, and users have the flexibility to include or remove  
 336 metrics according to their preference. The study includes 89 RCM simulations for precipitation,  
 337 and Fig. 5a shows the ranking of these models for different p norms. The lines corresponding  
 338 to each model give information about the model's ranking in different norms. The results  
 339 demonstrate that climate models are highly sensitive to p norms. Significant change in ranking  
 340 order is observed for the first four norms. Fig. 6 shows the percentage contribution of outliers  
 341 to the total error magnitude for models that have outliers. Median absolute deviation technique  
 342 (MAD) is used to identify outliers among the error metrics. Some of the models have only one  
 343 outlier (plots with a single solid line in Fig. 6) and other models have two outliers (plots with  
 344 both solid and dotted lines in Fig. 6). The percentage contribution of outliers increases as the p  
 345 norm increases, consistent with previous literature (Armstrong and Collopy, 1992; Hyndman  
 346 and Koehler, 2006). The study has used two parameters to indicate the capability of each norm  
 347 to differentiate between climate models - mean pairwise difference of the BM and the  
 348 difference between the maximum and minimum values of the BM. Figure 5b shows that both  
 349 parameters decrease as the p norm increases, indicating less differentiability. The results  
 350 suggest that the first norm (p=1) is the optimal norm to use as a metric in this study and will be  
 351 utilized in the following analyses.

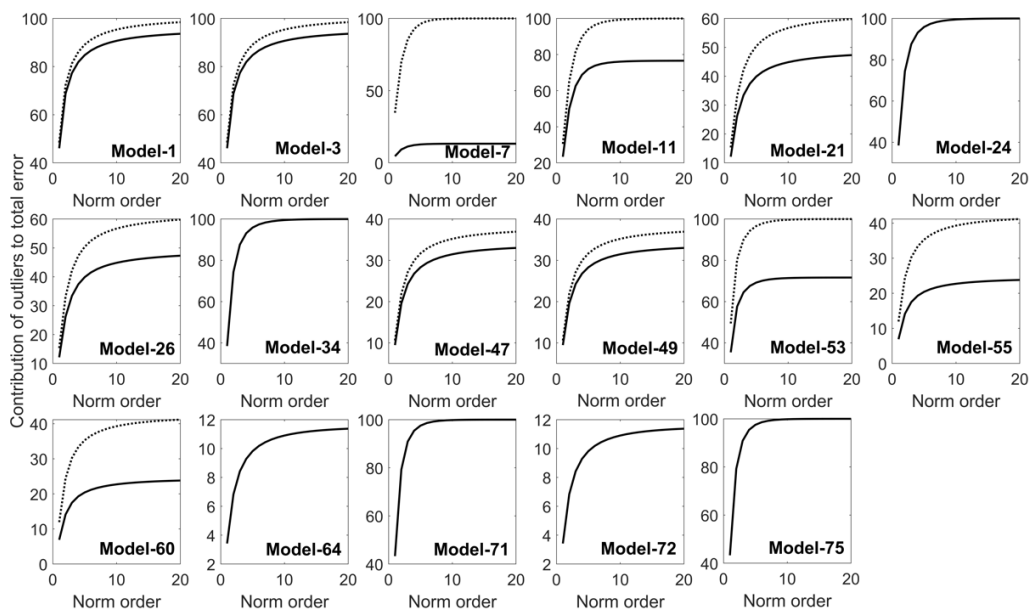
352 
$$\text{Bergen Metric (BM)} = \sqrt[p]{\begin{matrix} (0 - RMSE)^p + (0 - Bias)^p \\ + (1 - Standard\ deviation)^p \\ + (1 - Correlation\ coefficient)^p + (1 - Mean\ ratio)^p \end{matrix}} \quad (5)$$

353



354

355 **Figure 5:** a) The change in the ranking of the climate models with different norm order (p) b)  
 356 the change in the difference between the maximum and minimum distances and the average  
 357 distances with different norm order



358

359 **Figure 6:** The percentage contribution of outliers to the total error magnitude as a function of  
 360 norm order. The colours represent different outliers.

361 **4. Results**

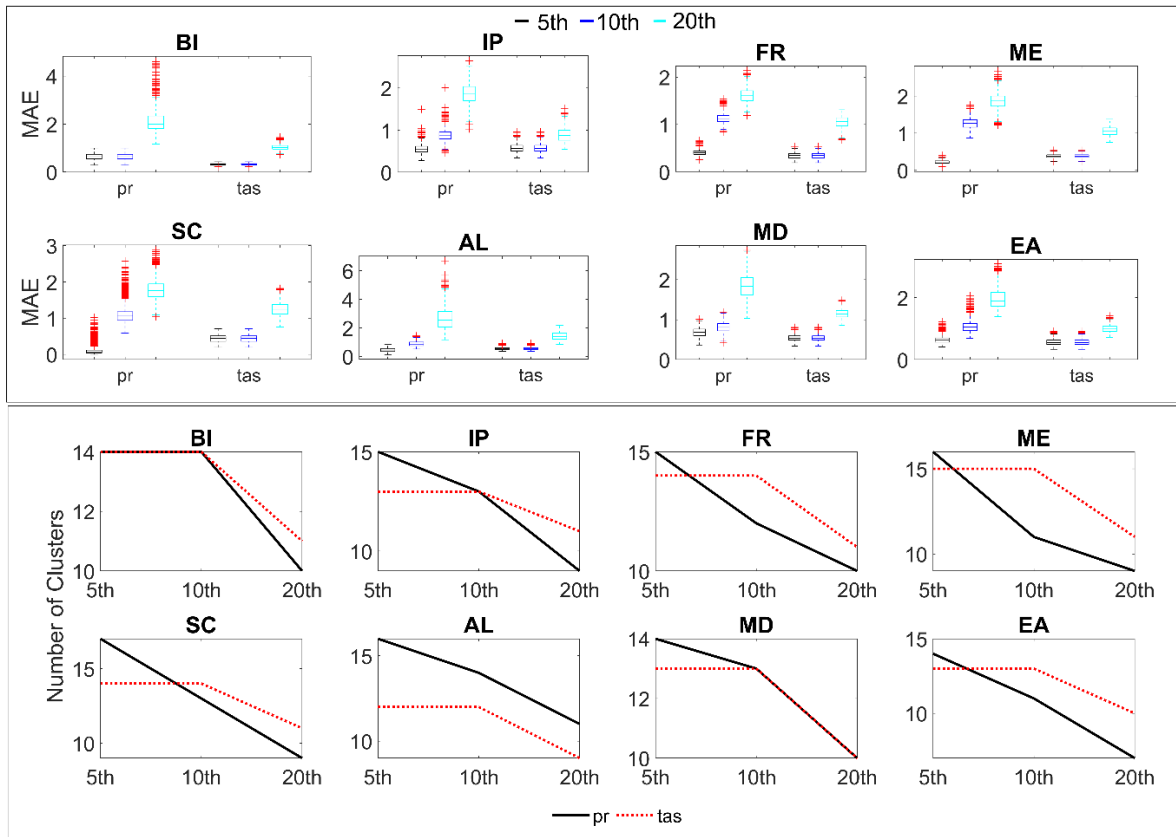
362 **4.1 Regional clustering of error metrics**

363 The study considers 38 error metrics (Table S3) which can take both positive and negative  
364 values as input. Similar to the models, the error metrics have been assigned a number (column  
365 1; Table S3) and the error metrics have been labelled as those numbers in some figures.

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

376 The 5th, 10th, and 20th percentiles are selected as potential thresholds to cluster the error  
377 metrics. However, users can select any number of thresholds for the sensitivity analysis. The  
378 clustering algorithm is allowed to run for all the thresholds to determine the optimal threshold.  
379 The efficiency of each cluster for a given threshold is represented by the mean of MAE over  
380 all the clusters. Another criterion used to determine the threshold is the number of clusters  
381 corresponding to each threshold. An increase in the percentile ( $q$ ) is expected to increase the  
382 MAE as the magnitude of threshold increases. Similarly, the number of clusters are expected  
383 to decrease as  $q$  increases as it can allow more error metrics into a cluster due to higher  
384 threshold magnitude. From Fig. 7, we conclude that the results are according to our  
385 expectations. It is found that increasing the percentile resulted in an increase in MAE and a  
386 decrease in the number of clusters. The 10th percentile is selected as the threshold to cluster  
387 the error metrics for both temperature and precipitation, as it has a smaller number of clusters  
388 compared to 5<sup>th</sup> percentile and less MAE compared to 20<sup>th</sup> percentile. The

389



390

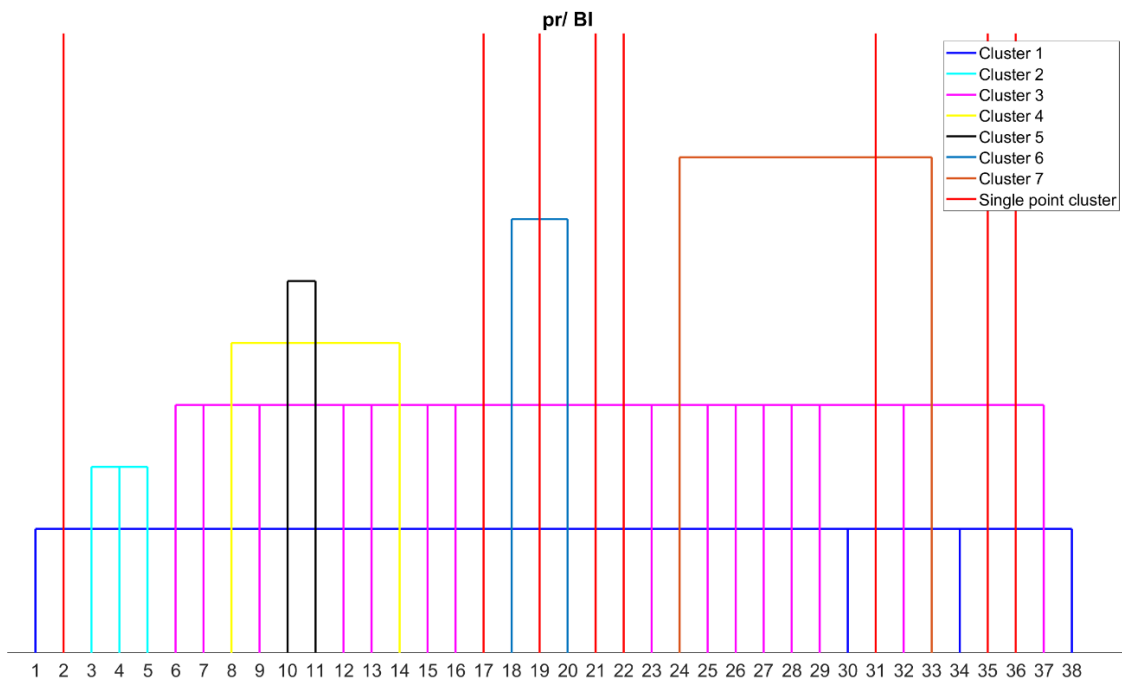
391 **Figure 7:** The variation in MAE (first box) and number of clusters (second box) corresponding  
 392 to 5<sup>th</sup>, 10<sup>th</sup> and 20<sup>th</sup> percentile for precipitation (pr) and temperature (tas) for all the eight regions

## 393 4.2 Results of clustering

### 394 4.2.1 Precipitation

395 For the British Isles region, the classification of 38 error metrics resulted in 15 clusters, with 8  
 396 error metrics being single point clusters due to their unique behaviour (Fig. 8). These 8 metrics  
 397 are d [2], (MB) R [17], MdE [19], MEE [21], MV [22], r2 [31], SGA [35], and R(Spearman)  
 398 [36]. The threshold for precipitation data is 6.35, indicating that all 8 error metrics produced  
 399 MAE values greater than 6.35 compared to the remaining 30 error metrics. RMSE [32] and its  
 400 variants such as normalized RMSE by IQR [25], mean [26] and range [27] are assigned to the  
 401 same cluster, as ED [7], IRMSE [9], MAE [13], MAPD [15], MASE [16], and MSE [23]. The  
 402 reason could be the L-norm framework which is used by most of the error metrics in this cluster.  
 403 D1 [3], d1 [4], and d(Mod.) [5] which share a similar framework, are also assigned to a single  
 404 cluster. Error metrics that evaluate the phase difference between observed and modelled data,  
 405 including ACC [1], R (Pearson) [30], SC [34], and M [38], are assigned to a single cluster.  
 406 H10(MAHE) [8] and MALE [14] share the same cluster as both metrics consider the difference  
 407 of logarithmic of the model and observed data to compute the error. Similarly, MDAE [18] and

408 MdSE [20] are assigned to a single cluster, as both metrics use the median of the difference  
 409 between observed and modelled data. However, MdE [19] is assigned to a different cluster as  
 410 it only considers the difference between observed and modelled data without bringing them to  
 411 the positive domain. NED [24] and SA [33] are found to be in the same cluster, as both metrics  
 412 are linearly associated while evaluating the model, even though their underlying frameworks  
 413 are somewhat different. Although ED [7] and NED [24] follow the L2 norm, they are not  
 414 assigned to the same cluster. This can be attributed to the normalisation of observed and  
 415 modelled data by their respective means in NED, as the statistical parameters such as mean is  
 416 sensitive to outliers, which can result in changes in ranking order.  
 417

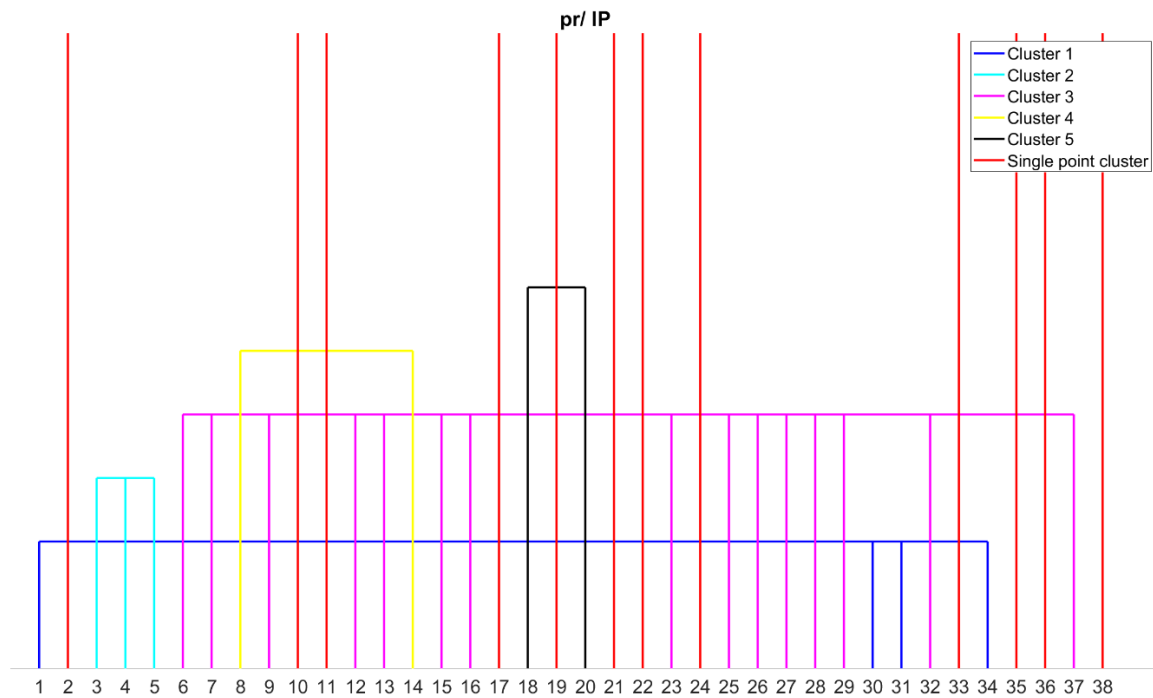


418  
 419 **Figure 8:** Clustering of error metrics using precipitation (pr) data for British Isles (BI) region.  
 420 Each error metric can be identified by the number using Table S3.

421 The Iberian Peninsula region is found to have 17 clusters, with 12 of them being single point  
 422 clusters (Fig. 9). Seven of the eight error metrics that are single point clusters in British Isles  
 423 are also single point clusters in Iberian Peninsula, except for r2 [31]. Five other error metrics:  
 424 NED [24], KGE (2009) [10], KGE (2012) [11], SA [33], and M [38] are also single point  
 425 clusters in Iberian Peninsula region. In British Isles, KGE (2009) [10] and KGE (2012) [11]  
 426 are assigned to the same cluster. The KGE (2012) is different from KGE (2009) since it used  
 427 the ratio of coefficient of variation between modelled and observed data instead of the ratio of  
 428 standard deviation to avoid the cross-correlation between bias and variability ratio. The



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



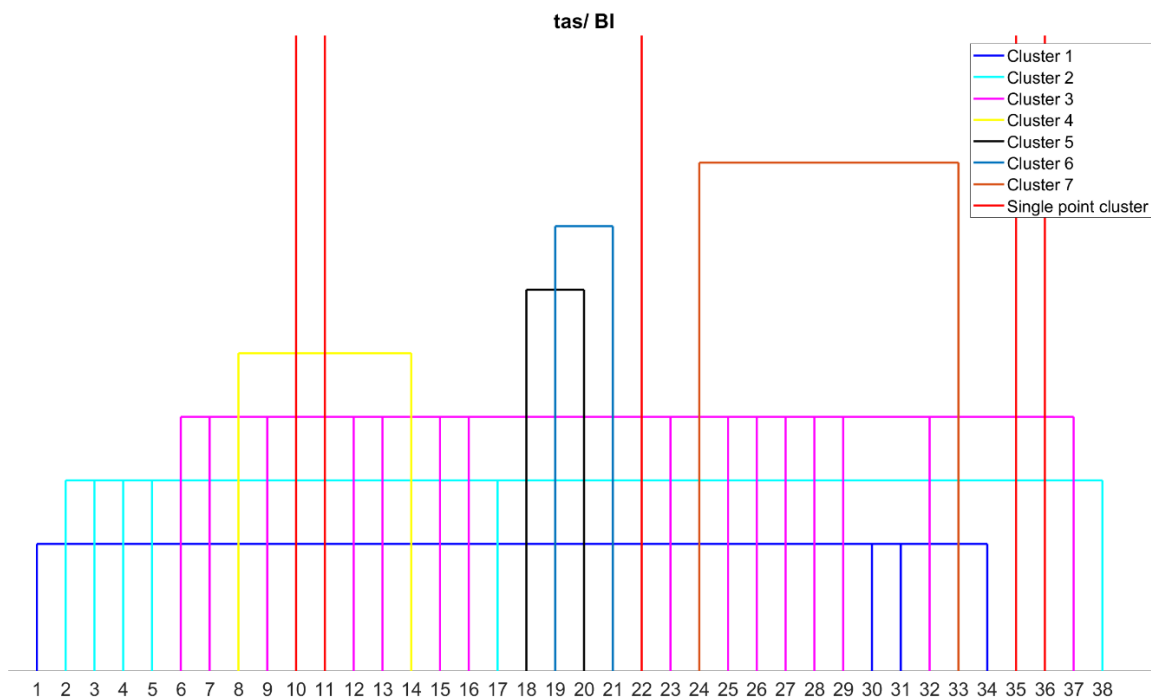
435  
 436 **Figure 9:** Clustering of error metrics using precipitation (pr) data for Iberian Peninsula (IP)  
 437 region. Each error metric can be identified by the number using Table S3.

438 As the results for the other 6 regions are similar to either the British Isles or the Iberian  
 439 Peninsula, we simply summarise their results here and refer the reader to the supplementary  
 440 material for further information. France (Fig. S2), Mid-Europe (Fig. S3), Scandinavia (Fig.  
 441 S4), Alps (Fig. S5), Mediterranean (Fig. S6) and Eastern Europe (Fig. S7) exhibit 15, 15, 16,  
 442 16, 17, and 14 clusters, respectively, with 8, 8, 10, 10, 12, and 6 single point clusters. France  
 443 and Mid-Europe have the same clusters as the British Isles, and the Mediterranean has the same  
 444 clusters as Iberian Peninsula. Scandinavia has clusters similar to British Isles, except that M  
 445 [38] is a single point cluster and  $r^2$  [31] has been assigned to the correlation metrics cluster in  
 446 Scandinavia. The Alps also has clusters similar to British Isles, except KGE (2009) [10] and  
 447 KGE (2012) [11] are single point clusters. Eastern Europe also has clusters similar to British

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

#### 450 4.2.2 Temperature

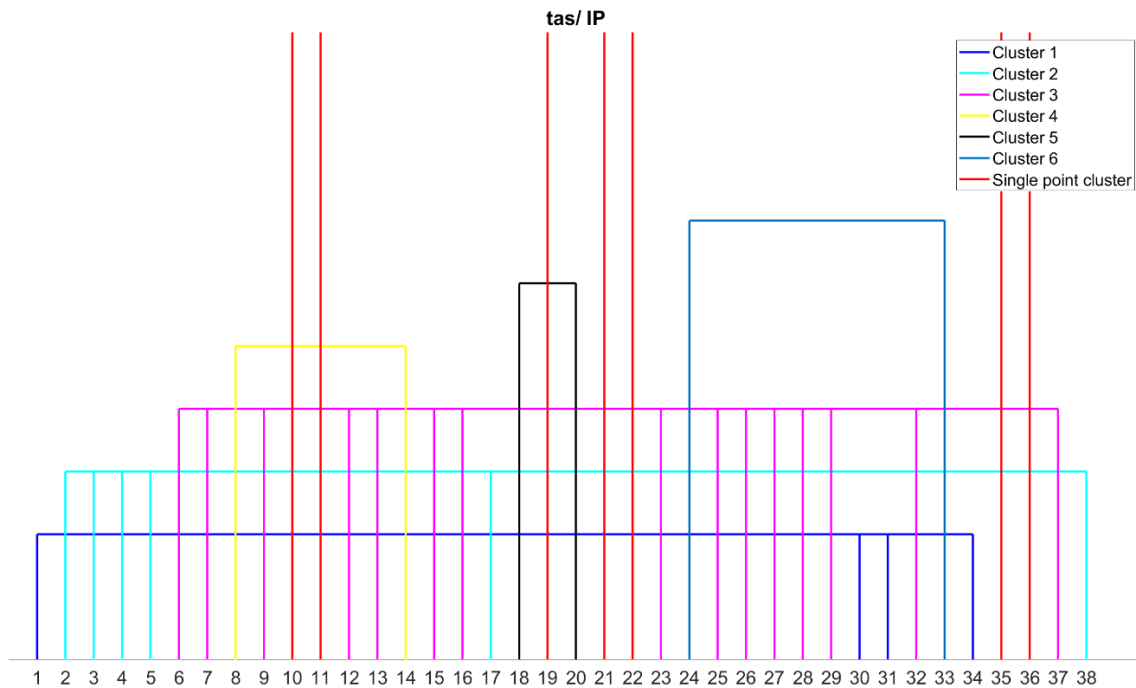
451 Compared to precipitation data, temperature data has a lower number of clusters, which can be  
 452 attributed to the lower variability in temperature data. The clustering of error metrics for British  
 453 Isles is shown in Fig. 10. For British Isles, 12 clusters are identified, with 5 single point clusters,  
 454 namely KGE(2009) [10], KGE(2012) [11], MV [22], SGA [35], and R(Spearman) [36]. Similar  
 455 to precipitation clusters, several error metrics, including ED [7], IRMSE [9], MAE [13], MAPD  
 456 [15], MASE [16], MSE [23], NRMSE(IQR) [25], NRMSE(mean) [26], NRMSE(range) [27]  
 457 and RMSE [32] are assigned to the same cluster.



458  
 459 **Figure 10:** Clustering of error metrics using temperature (tas) data for British Isles (BI) region.  
 460 Each error metric can be identified by the number using Table S3.

461 The correlation metrics, such as ACC [1], r2 [31], SCO [34], and R(Pearson) [36] belong to  
 462 the same cluster. France (Fig. S8) and Mid-Europe (Fig. S9) have the same cluster as British  
 463 Isles for temperature data. For Iberian Peninsula (Fig. 11), 13 different clusters are identified,  
 464 with 7 single point clusters, including MdE [19] and MEE [21] in addition to the 5 single point  
 465 clusters from British Isles. The remaining clusters are similar to those in British Isles.  
 466 Mediterranean (Fig. S10) has the same cluster as Iberian Peninsula for temperature data, with  
 467 13 clusters and 7 single point clusters. Scandinavia (Fig. S11) and Eastern Europe (Fig. S12)

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



471  
 472 **Figure 11:** Clustering of error metrics using temperature (tas) data for Iberian Peninsula (IP)  
 473 region. Each error metric can be identified by the number using Table S3.

474 **4.3 Bergen Metrics**

475 A Bergen metric is computed for all eight regions using the respective clusters for both  
 476 precipitation and temperature. A single metric is chosen from each cluster randomly; Random  
 477 selection demonstrated no discernible impact on the ranking (see Text S2). Although computed  
 478 for all 89 regional climate models, this paper focuses on discussing only one climate model for  
 479 both precipitation and temperature. The CLM Community (CLMCom) regional model from  
 480 ICHEC-EC-EARTH for r3i1p1 realisation is discussed as it performed best at over 25 grid  
 481 points in 5 regions and more than 2 grid points in seven regions. For the temperature variable,  
 482 the CLMCom model form CCCma-CanESM2 model for r1i1p1 realisation is discussed, as it  
 483 performed best at over 25 grid points in seven regions.

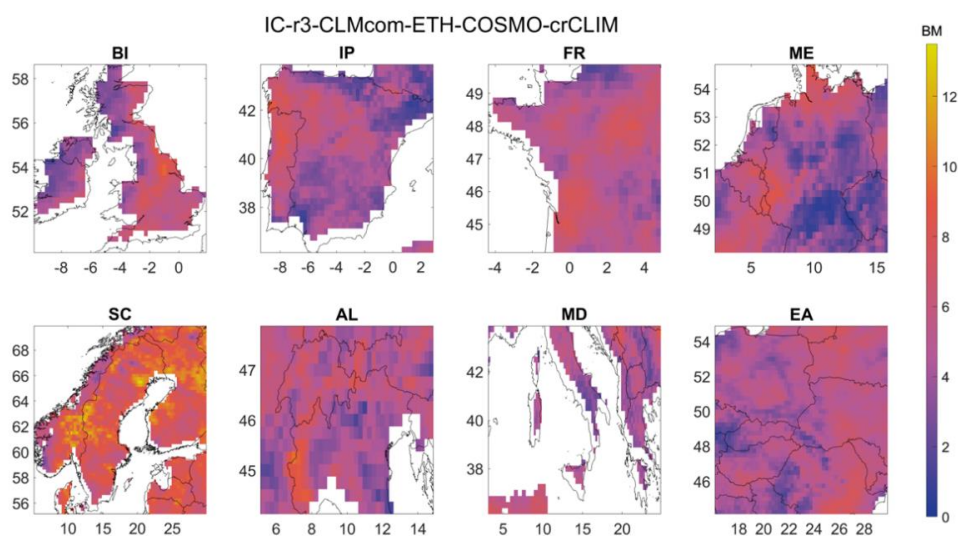
484 **4.3.1 Precipitation**

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

488 MdAE, NED, d, MB(R), MdE, MEE, MV, r2, SGA, and R(Spearman). Figure 12 provides an  
 489 overview of the spatial distribution of the BM for all eight regions, while the spatial distribution  
 490 of each of these metrics is shown in Fig. 13 for the British Isles region.

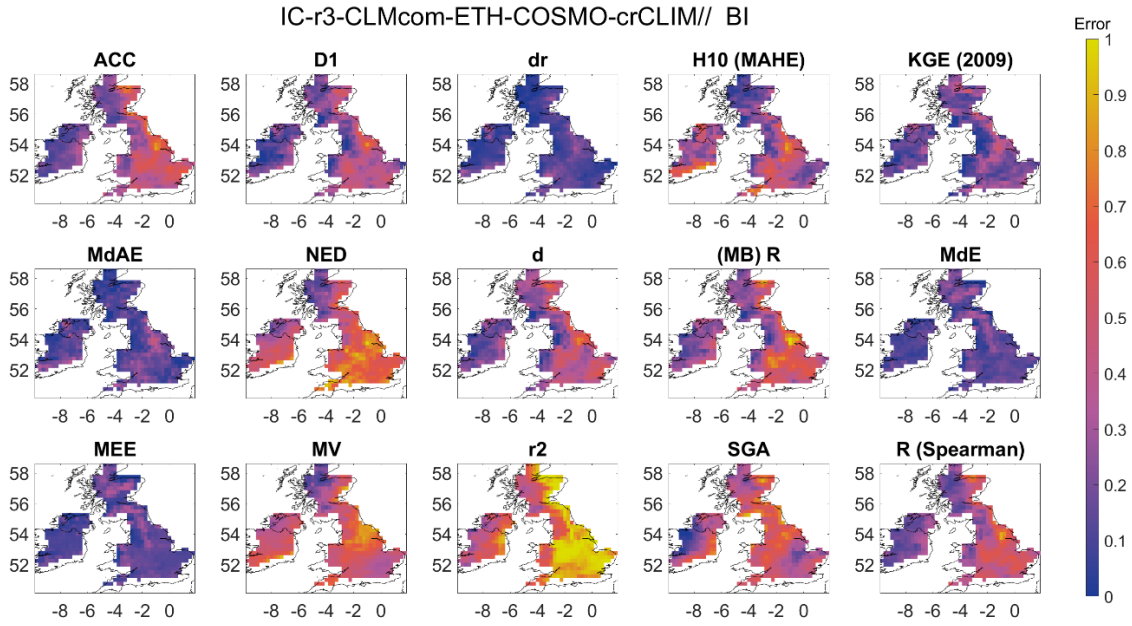
491 The magnitude of BM ranges from 0 to 13, with a score of 0 indicating good performance by  
 492 the model. Based on the results, the CLMCom model performed well in the western part of  
 493 British Isles, as indicated by the BM. This is a result of the good performance of most of the  
 494 individual metrics that comprise the Bergen Metric. This is shown in Fig. 13. There are some  
 495 contradictory results from different error metrics in the eastern region. While all 13 metrics  
 496 indicate good performance, the MV, r2 and NED indicate very bad performance by the model.  
 497 The use of individual error metrics can provide meaningful insights into the performance of  
 498 the model in different regions. For example, metrics such as dr, MdAE, MdE, and MEE  
 499 indicate good performance in the southeastern region, while R(Spearman) indicates bad  
 500 performance by the CLMCom model which implies that the phase difference is significant  
 501 between observed and modelled data in this region. It is worth noting that some metrics, such  
 502 as r2 and R(Spearman), may provide different results even though they share a similar  
 503 framework. R(Spearman) only tells how well the modelled data follow the observed data while  
 504 r2 indicate how well the data represents the line of best fit (<https://tinyurl.com/y52r3xed>;  
 505 <https://tinyurl.com/yk2jmsxt>). Overall, the use of multiple error metrics and the analysis of  
 506 individual metrics can provide a more comprehensive assessment of the model's performance,  
 507 particularly in regions where different metrics provide conflicting results.

508



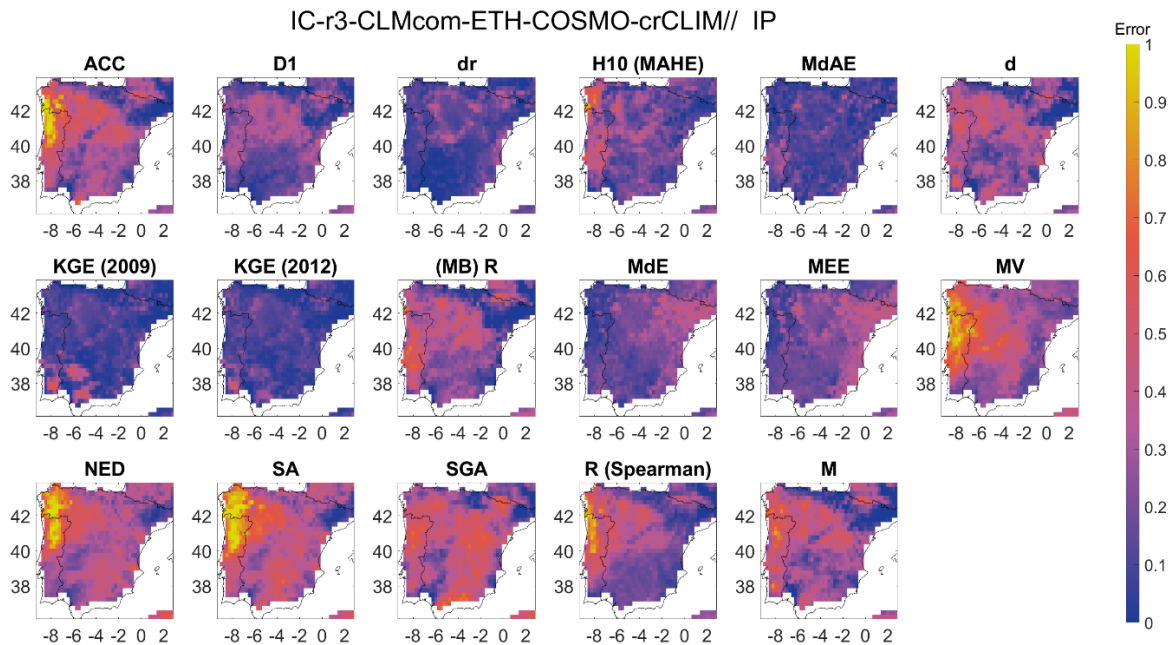
509

510 **Figure 12:** Spatial distribution of Bergen metric using precipitation data for all the eight  
 511 regions



512

513 **Figure 13:** Spatial distribution of the error metrics used to compute the Bergen metric for  
 514 precipitation and for British Isles (BI) region. The error metrics have been labelled by the  
 515 abbreviation and the corresponding error metrics can be identified from Table S3.



516

517 **Figure 14:** Spatial distribution of the error metrics used to compute the Bergen metric for  
 518 precipitation and for Iberian Peninsula (IP) region. The error metrics have been labelled by the  
 519 abbreviation and the corresponding error metrics can be identified from Table S3.

520 Figure 14 shows a Bergen metric for Iberian Peninsula applied to the CLMCom model, which  
 521 is based on 17 error metrics obtained from each cluster. These metrics, including ACC, D1, dr,

522 H10 (MAHE), MdAE, d, KGE (2009), KGE (2012), MB (R), MdE, MEE, MV, NED, SA,  
523 SGA, R (Spearman) and M, are presented in Fig. 14. The results indicate that the model  
524 performs relatively better in the northeast and southeast regions compared to the western region  
525 (see Fig. 12), possibly due to the influence of certain metrics such as ACC, R (Spearman), MV,  
526 NED, and SA. Additionally, while KGE (2009) and KGE (2012) exhibit similar spatial error  
527 patterns, further analysis in the southern region reveals the differences in the magnitude of  
528 error. Interestingly, despite their similarity, KGE (2009) and KGE (2012) are classified into  
529 different clusters based on a threshold MAE of 5.41, used to determine cluster membership.

530

531 France (Fig. S14), and Mid-Europe (Fig. S15) have the same clusters as the British Isles, and  
532 therefore the same error metrics used in British Isles are used to calculate the Bergen metric  
533 for France and Mid-Europe. The Bergen metric indicates an average performance of the model  
534 for the entire study region of France (see Fig. 12). While  $r^2$  shows a very poor performance of  
535 the model for France, MEE metric shows a completely opposite trend, indicating a very good  
536 performance of the model. Similar disagreement between  $r^2$  and MEE is also observed in the  
537 British Isles. On the other hand, SGA, which compares the shape of the two signals, shows an  
538 average performance by the model. In terms of the spatial distribution of error, the Bergen  
539 metric shows lower error magnitudes for MEE in the southeast part of the study region.

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

545 Fig. S16 and Fig. 12 suggest that the CLMCom model does not perform well for Scandinavia.  
546 However, some error metrics, including dr, MdAE, MdE, and MEE, show good performance  
547 in the southern part of the region. Although MdAE, MdE, and MEE are assigned to different  
548 clusters, they exhibit similar spatial distributions of error. It is worth noting that despite the  
549 similarity, the three error metrics are in different clusters due to their higher MAE between  
550 them. For the Alps, the Bergen metric indicates a relatively good performance of the CLMCom  
551 model. It can be observed in Fig. S17, all metrics except  $r^2$  show good performance for the  
552 model.

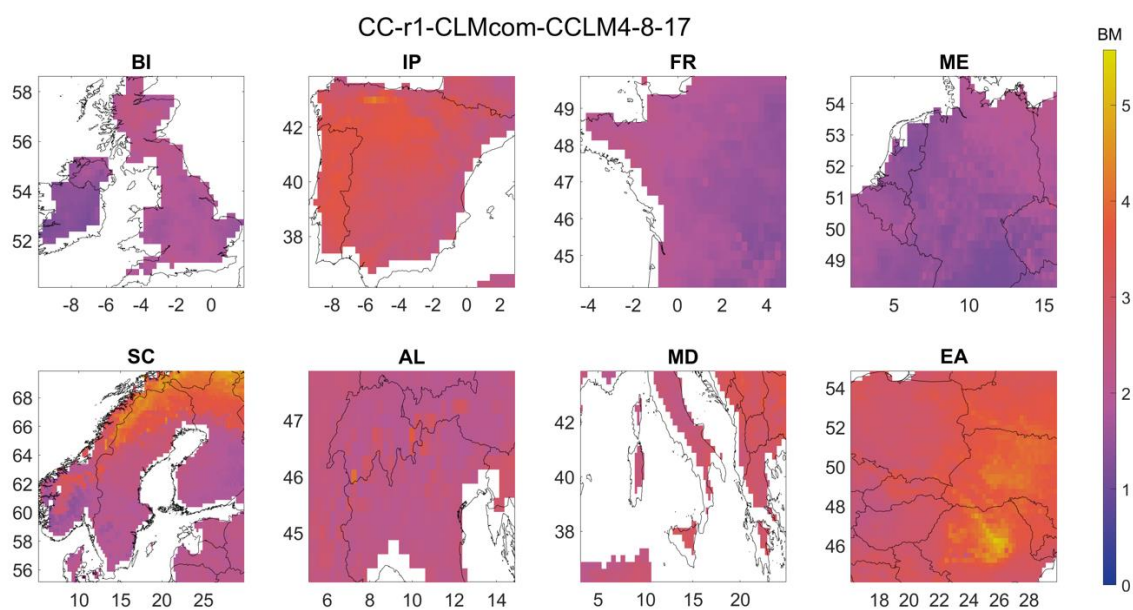
553 The Mediterranean has the same clusters as the Iberian Peninsula, and the spatial distribution  
554 of each metric for the Mediterranean is presented in Fig. S18. The Bergen metric for the  
555 CLMCom model suggests an average performance for the entire Mediterranean region. Some

556 of the error metrics, such as KGE (2009), KGE (2012), dr, and MdAE, indicate good model  
 557 performance. However, metrics such as SGA, SA, and NED, show relatively poor performance  
 558 of the model.

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

### 566 4.3.2 Temperature

567 For temperature, we focus on the CLM Community (CLMCom) regional model driven by  
 568 ICHEC-EC-EARTH to demonstrate the application of Bergen metrics for temperature. The  
 569 spatial distribution of BM is shown in Fig. 15, which indicates average performance by the  
 570 model, except in certain areas like northern part of Scandinavia, central part of Eastern Europe  
 571 and western part of Iberian Peninsula, where the performance is bad. The British Isles (Fig.  
 572 16), France (Fig. S20), and Mid-Europe (Fig. S21) regions have 12 clusters, and 12 error  
 573 metrics, including ACC, d, dr, H10(MAHE), MdAE, MdE, NED, KGE(2009), KGE(2012),  
 574 MV, SGA, and R(Spearman) are used to compute the Bergen metric for these regions.

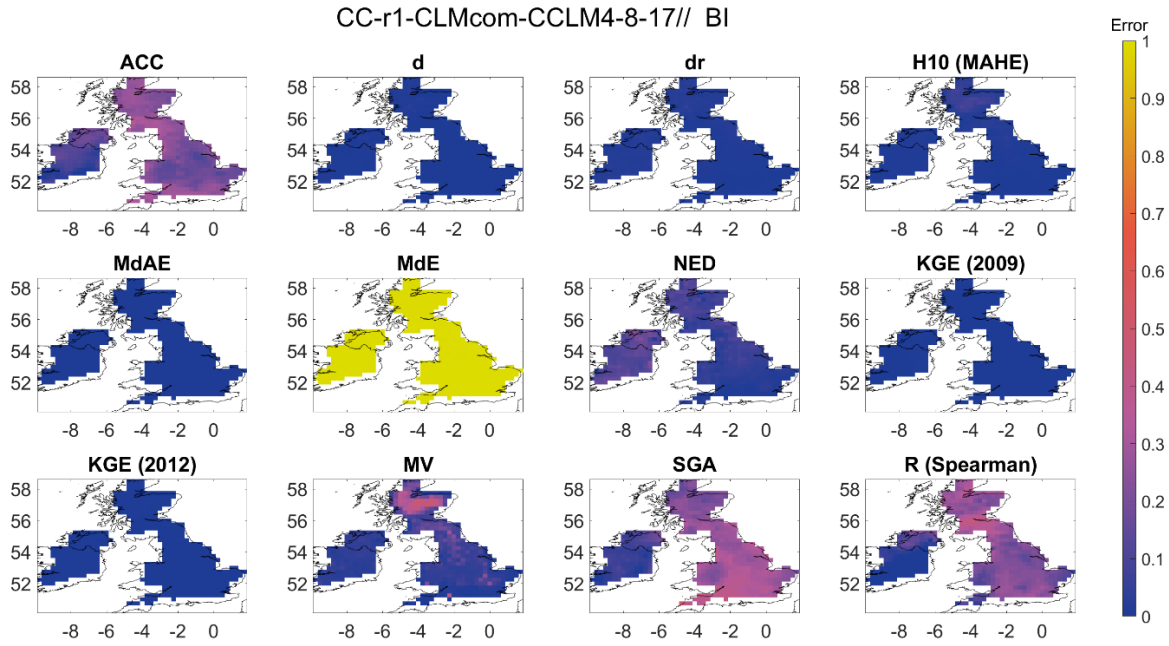


575  
 576 **Figure 15:** Spatial distribution of Bergen metric using temperature data for all the eight regions

577 The Scandinavia (Fig. S22) and Eastern Europe (Fig. S23) regions have 14 clusters and all the  
578 error metrics from British Isles, along with VE and SA, are used to compute the Bergen metric  
579 for these regions. The Iberian Peninsula (Fig. 17) and Mediterranean (Fig. S24) regions have  
580 the same cluster, with a total of 13 clusters and all the error metrics from British Isles, plus  
581 MEE, are used to compute the Bergen metric. The Alps (Fig. S25) region has 15 clusters, with  
582 all the error metrics from Scandinavia, including MEE, used to compute the Bergen metric.  
583 MdE and MEE consistently indicate very bad model performance for all the regions, while the  
584 other metrics indicate relatively good performance. This suggests that the mean and median of  
585 the modelled data tend to underestimate/overestimate the observed mean and median,  
586 respectively. Histograms in Fig. 18 further investigate this, showing that the error values for  
587 ACC are more evenly distributed in the Iberian Peninsula region and close to its ideal point 1,  
588 while the source errors for MdE and MEE are concentrated between -0.5 to -1.5, resulting in  
589 most of the error values being concentrated between 0.9 to 1 after normalization. The source  
590 error represents the distance between the ideal values and actual magnitude after normalization.  
591 Similar patterns can be observed in the other regions for temperature.

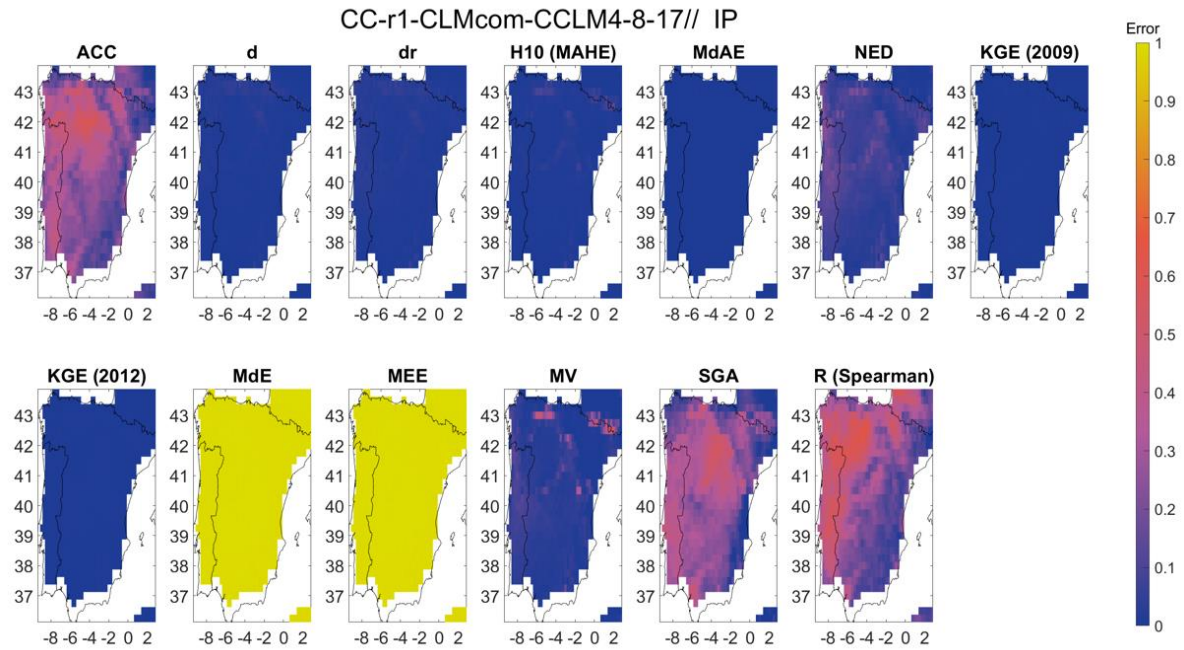
592 To illustrate inter-model variability, a random grid point (50.125, 1.875) is selected. The  
593 Bergen metric is calculated for both precipitation and temperature at this grid point, and models  
594 are ranked based on the Bergen metric (Fig. 19). The Bergen metric ranges from 2.29 to 11.39  
595 for precipitation and 1.85 to 8.37 for temperature. Notably, with a Bergen metric value of 2.29,  
596 ETH-COSMO (Model 6) is identified as performing well for precipitation. Similarly, with a  
597 Bergen metric value of 2.29, GERICS-REMO2015 (Model 16) is recognized for its good  
598 performance in temperature. The proposed metric offers a valuable tool for assessing the  
599 performance of climate models.





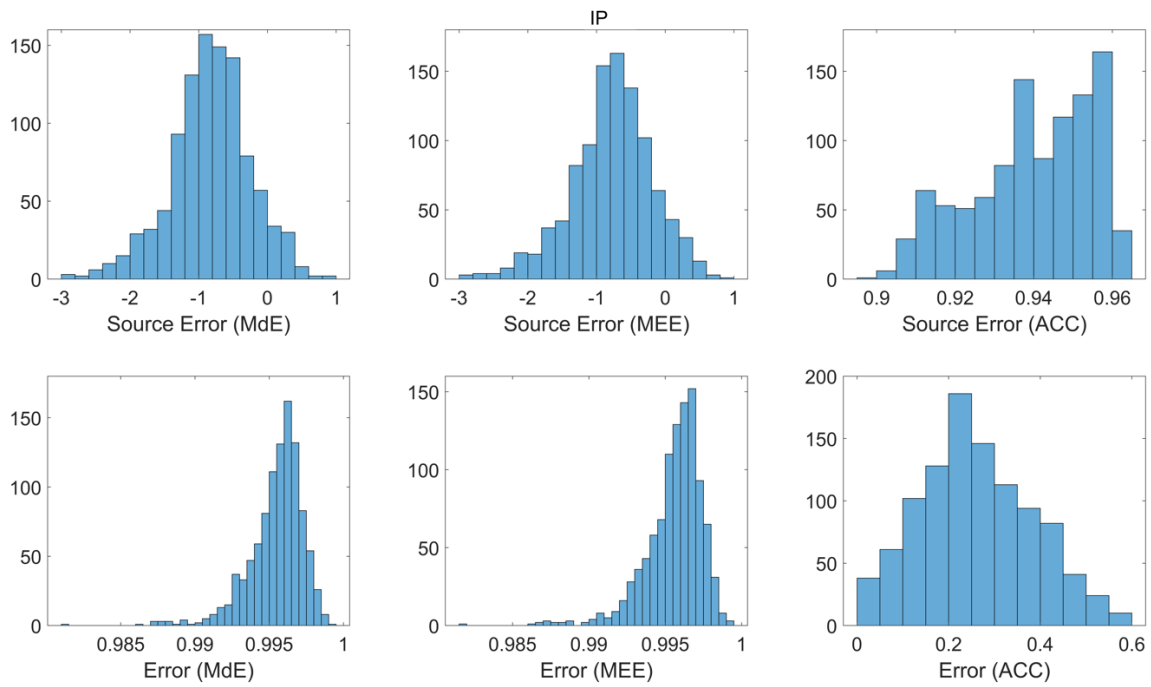
600

601 **Figure 16:** Spatial distribution of the error metrics used to compute the Bergen metric for  
 602 temperature and for British Isles (BI) region. The error metrics have been labelled by the  
 603 abbreviation and the corresponding error metrics can be identified from Table S3.



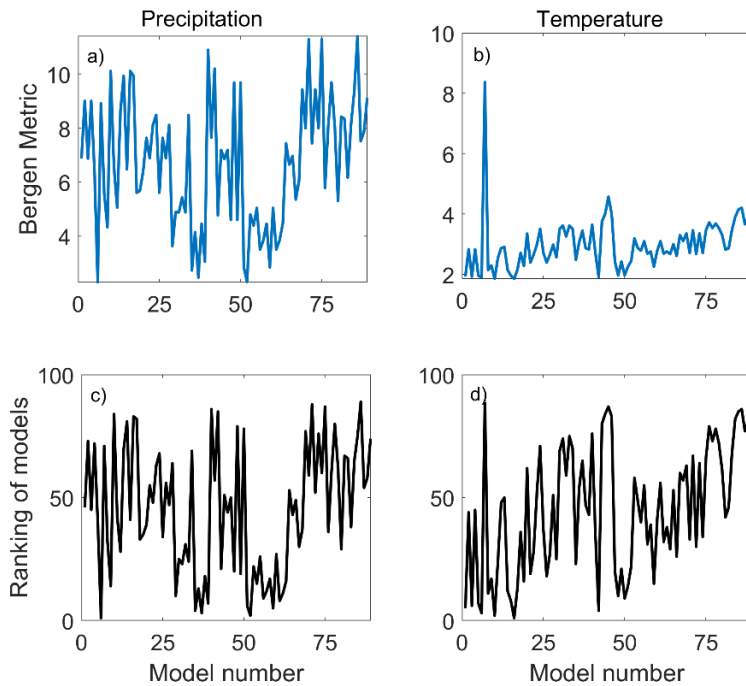
604

605 **Figure 17:** Spatial distribution of the error metrics used to compute the Bergen metric for  
 606 temperature and for Iberian Peninsula (IP) region. The error metrics have been labelled by the  
 607 abbreviation and the corresponding error metrics can be identified from Table S3.



608

609 **Figure 18:** Histogram plot of error and source error for MdE, MEE and ACC for Iberian  
 610 Peninsula region (IP).



611

612 **Figure 19:** The Bergen metric for precipitation (a) and temperature (b) for all 89 climate  
 613 models, along with the ranking of each model based on the Bergen metric for precipitation (c)  
 614 and temperature (d), at a grid point (50.125, 1.875).

## 5. Conclusions

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

The study also shows that the number of error metrics used in Bergen Metrics can be reduced using a non-parametric clustering technique. Although several clustering techniques are already available in the literature, they come with certain requirements. Either they require the number of clusters before running the algorithm or information on the class label of the feature vector. The adopted clustering technique tries to identify the natural cluster present in the data. The mean absolute error based on ranking order is used as a dissimilarity index to assign error metrics to different clusters. The technique also has a threshold parameter  $5^{\text{th}}$ ,  $10^{\text{th}}$  and  $20^{\text{th}}$  are selected as candidates for threshold parameter and  $10^{\text{th}}$  percentile of the D matrix is adopted as a threshold in this study. It is selected because increase in threshold ( $20^{\text{th}}$  percentile) resulted in increase in MAE and decrease in number of clusters, whereas, decrease in threshold ( $5^{\text{th}}$  percentile) resulted in decrease in MAE and increase in number of clusters and the study chose a middle ground. However, users can investigate different values of q before choosing the threshold. The clustering technique is compared with the K-means clustering approach and it is found that the non-parametric technique has lower MAE compared to the K-means approach. The clustering is performed for all the eight regions and those are British Isles, Iberian Peninsula, France, Mid-Europe, Scandinavia, Alps, Mediterranean and Eastern Europe. For precipitation, 15, 17, 15, 15, 16, 15, 17, and 14 clusters are obtained for the eight regions, respectively. For temperature, 12, 13, 12, 12, 14, 15, 13, and 14 clusters are obtained for the eight regions, respectively.

A single error metric from each cluster can be chosen randomly as a component to be used in the calculation of a Bergen Metric. We have shown that random selection does not have any effect on the ranking order produced by a Bergen Metric. The Bergen Metric which uses the L1 framework is found to be less sensitive to outliers compared to the other norms and more stable in higher dimensional space. Bergen Metrics are a multivariate error functions that can take any number of error metrics of different variables as shown in the last section. It can be further modified for a weighting-based metric that can allow the user to give more weightage to particular metrics depending on the requirement of the study. While some metrics show good

649 performance in certain regions, others indicate poor performance. It is also important to observe  
650 how a single metric can influence and change the ranking of climate models. Bergen metrics  
651 provide a comprehensive evaluation of the model's performance, which is useful for identifying  
652 the strengths and weaknesses of the model in different contexts. It is also crucial to underscore  
653 that our proposed metric evaluates the magnitude differences between modeled and reference  
654 data, prioritizing this aspect over spatial and temporal patterns. The application of this metric  
655 should be approached with careful consideration.

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

663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682

683 **Data and Code availability**

684 The EURO-CORDEX data used in this work are obtained from the Earth System Grid  
685 Federation server. The reference precipitation and temperature data is available at  
686 [https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-pressure-levels-monthly-](https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-pressure-levels-monthly-means-preliminary-back-extension?tab=form)  
687 [means-preliminary-back-extension?tab=form](https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-pressure-levels-monthly-means-preliminary-back-extension?tab=form)

688 The code for clustering the error metrics is available at  
689 <https://doi.org/10.5281/zenodo.10518064>

690

691 **Author contributions**

692 AS developed the methodology and performed the formal analysis. PM supervised the research  
693 activity planning and execution. AS prepared the first draft of manuscript. All authors  
694 contributed to editing and reviewing the manuscript.

695

696 **Competing interests**

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

698

699 **Acknowledgements**

700 The FRONTIER project has received funding from the Research Council of Norway (project  
701 number 301777). We thank James Done and Andreas Prein for their advice and critical  
702 comments regarding the work.

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717 **References**

- 718 Aggarwal, C. C., Hinneburg, A., & Keim, D. A.: On the surprising behavior of distance metrics in high  
719 dimensional space. In *International conference on database theory* (pp. 420-434). Springer, Berlin, Heidelberg,  
720 DOI: 10.1007/3-540-44503-X\_27, 2001.
- 721 Ahmed, K., Sachindra, D. A., Shahid, S., Demirel, M. C., & Chung, E. S.: Selection of multi-model ensemble of  
722 general circulation models for the simulation of precipitation and maximum and minimum temperature based on  
723 spatial assessment metrics. *Hydrology and Earth System Sciences*, 23(11), 4803-4824,  
724 <https://doi.org/10.5194/hess-23-4803-2019>, 2019.
- 725 Armstrong, J. S., & Collopy, F.: Error measures for generalizing about forecasting methods: Empirical  
726 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),  
727 1992.
- 728 Baker, N. C., & Taylor, P. C.: A framework for evaluating climate model performance metrics. *Journal of*  
729 *Climate*, 29(5), 1773-1782, <https://doi.org/10.1175/JCLI-D-15-0114.1>, 2016
- 730 Bellomo, K., Angeloni, M., Corti, S., & von Hardenberg, J.: Future climate change shaped by inter-model  
731 differences in Atlantic meridional overturning circulation response. *Nature Communications*, 12(1), 1-10,  
732 <https://doi.org/10.1038/s41467-021-24015-w>, 2021.
- 733 Benestad, R. E., Mezghani, A., Lutz, J., Dobler, A., Parding, K. M., & Landgren, O. A.: Various ways of using  
734 empirical orthogonal functions for climate model evaluation. *Geoscientific Model Development*, 16(10), 2899-  
735 2913, <https://doi.org/10.5194/gmd-16-2899-2023>, 2023.
- 736 Boberg, F., Berg, P., Thejll, P. et al. Improved confidence in climate change projections of precipitation further  
737 evaluated using daily statistics from ENSEMBLES models. *Clim Dyn*, 35, 1509–1520,  
738 <https://doi.org/10.1007/s00382-009-0683-8>, 2010
- 739 Boberg, F., Berg, P., Thejll, P. et al. Improved confidence in climate change projections of precipitation evaluated  
740 using daily statistics from the PRUDENCE ensemble. *Clim Dyn*, 32, 1097–1106, <https://doi.org/10.1007/s00382-008-0446-y>, 2009
- 742 Chai, T., & Draxler, R. R.: Root mean square error (RMSE) or mean absolute error (MAE)?—Arguments against  
743 avoiding RMSE in the literature. *Geoscientific model development*, 7(3), 1247-1250,  
744 <https://doi.org/10.5194/gmd-7-1247-2014>, 2014.
- 745 Christensen, J. H., & Christensen, O. B.: A summary of the PRUDENCE model projections of changes in  
746 European climate by the end of this century. *Climatic change*, 81(Suppl 1), 7-30, <https://doi.org/10.1007/s10584-006-9210-7>, 2007.
- 748 Collier, N., Hoffman, F. M., Lawrence, D. M., Keppel-Aleks, G., Koven, C. D., Riley, W. J., ... & Randerson, J.  
749 T.: The International Land Model Benchmarking (ILAMB) system: design, theory, and implementation. *Journal*  
750 *of Advances in Modeling Earth Systems*, 10(11), 2731-2754, <https://doi.org/10.1029/2018MS001354>, 2018.
- 751 Dembélé, M., Hrachowitz, M., Savenije, H. H., Mariéthoz, G., & Schaeffli, B.: Improving the predictive skill of a  
752 distributed hydrological model by calibration on spatial patterns with multiple satellite data sets. *Water resources*  
753 *research*, 56(1), e2019WR026085, <https://doi.org/10.1029/2019WR026085>, 2020.
- 754 Demirel, M. C., Mai, J., Mendiguren, G., Koch, J., Samaniego, L., & Stisen, S.: Combining satellite data and  
755 appropriate objective functions for improved spatial pattern performance of a distributed hydrologic model.  
756 *Hydrology and Earth System Sciences*, 22(2), 1299-1315, <https://doi.org/10.5194/hess-2017-570>, 2018.
- 757 Flato, G., Marotzke, J., Abiodun, B., Braconnot, P., Chou, S. C., Collins, W., Cox, P., Driouech, F., Emori, S.,  
758 Eyring, V., Forest, C., Gleckler, P., Guilyardi, E., Jakob, C., Kattsov, V., Reason, C., & Rummukainen, M.:  
759 Evaluation of climate models. In *Climate Change 2013: the physical science basis. Contribution of Working*  
760 *Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*, 741-866, [Stocker,

- 761 T.F., D. Qin, G.-K. Plattner, M. Tignor, S.K. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex and P.M. Midgley  
762 (eds.].Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA, 2013.
- 763 Giot, Olivier, Piet Termonia, Daan Degrauwe, Rozemien De Troch, Steven Caluwaerts, Geert Smet, Julie  
764 Berckmans et al.: Validation of the ALARO-0 model within the EURO-CORDEX framework. *Geoscientific*  
765 *Model Development* , 1143-1152. doi:10.5194/gmd-9-1143-2016, 2016.
- 766 Graham, R. M., Cohen, L., Ritzhaupt, N., Segger, B., Graversen, R. G., Rinke, A., Walden, V.P., Granskog, M.A.,  
767 & Hudson, S. R.: Evaluation of six atmospheric reanalyses over Arctic sea ice from winter to early  
768 summer. *Journal of Climate*, 32(14), 4121-4143, <https://doi.org/10.1175/JCLI-D-18-0643.1>, 2019.
- 769 Gupta, H. V., Kling, H., Yilmaz, K. K., & Martinez, G. F.: Decomposition of the mean squared error and NSE  
770 performance criteria: Implications for improving hydrological modelling. *Journal of hydrology*, 377(1-2), 80-91,  
771 <https://doi.org/10.1016/j.jhydrol.2009.08.003>, 2009.
- 772 Hartigan, J. A., & Wong, M. A.: Algorithm AS 136: A k-means clustering algorithm. *Journal of the royal*  
773 *statistical society. series c (applied statistics)*, 28(1), 100-108, <https://doi.org/10.2307/2346830>, 1979.
- 774 He, X., Lei, X. D., & Dong, L. H.: How large is the difference in large-scale forest biomass estimations based on  
775 new climate-modified stand biomass models?. *Ecological Indicators*, 126, 107569,  
776 <https://doi.org/10.1016/j.ecolind.2021.107569>, 2021.
- 777 Hu, Z., Chen, X., Zhou, Q., Chen, D., & Li, J.: DISO: A rethink of Taylor diagram. *International Journal of*  
778 *Climatology*, 39(5), 2825-2832, <https://doi.org/10.1002/joc.5972>, 2019.
- 779 Hyndman, R. J., & Koehler, A. B.: Another look at measures of forecast accuracy. *International journal of*  
780 *forecasting*, 22(4), 679-688, <https://doi.org/10.1016/j.ijforecast.2006.03.001>, 2006.
- 781 IPCC: Climate Change 2021: The Physical Science Basis. Contribution of Working Group I to the Sixth  
782 Assessment Report of the Intergovernmental Panel on Climate Change, edited by: Masson-Delmotte, V., Zhai, P.,  
783 Pirani, A., Connors, S. L., Péan, C., Berger, S., Caud, N., Chen, Y., Goldfarb, L., Gomis, M. I., Huang, M.,  
784 Leitzell, K., Lonnoy, E., Matthews, J. B. R., Maycock, T. K., Waterfield, T., Yelekçi, O., Yu, R., & Zhou, B.,  
785 Cambridge University Press, [https://report.ipcc.ch/ar6/wg1/IPCC\\_AR6\\_WGI\\_FullReport.pdf](https://report.ipcc.ch/ar6/wg1/IPCC_AR6_WGI_FullReport.pdf), 2021a.
- 786 IPCC: Summary for Policymakers, in: Climate Change 2021: The Physical Science Basis. Contribution of  
787 Working Group I to the Sixth Assessment Report of the Intergovernmental Panel on Climate Change, edited by:  
788 Masson-Delmotte, V., Zhai, P., Pirani, A., Connors, S. L., Péan, C., Berger, S., Caud, N., Chen, Y., Goldfarb, L.,  
789 Gomis, M. I., Huang, M., Leitzell, K., Lonnoy, E., Matthews, J. B. R., Maycock, T. K., Waterfield, T., Yelekçi,  
790 O., Yu, R., & Zhou, B., Cambridge University Press,  
791 [https://www.ipcc.ch/report/ar6/wg1/downloads/report/IPCC\\_AR6\\_WGI\\_SPM.pdf](https://www.ipcc.ch/report/ar6/wg1/downloads/report/IPCC_AR6_WGI_SPM.pdf), 2021b.
- 792 Jackson, E. K., Roberts, W., Nelsen, B., Williams, G. P., Nelson, E. J., & Ames, D. P.: Introductory overview:  
793 Error metrics for hydrologic modelling—A review of common practices and an open source library to facilitate use  
794 and adoption. *Environmental Modelling & Software*, 119, 32-48, <https://doi.org/10.1016/j.envsoft.2019.05.001>,  
795 2019.
- 796 Kalmár, T., Pieczka, I., & Pongrácz, R.: A sensitivity analysis of the different setups of the RegCM4.5 model for  
797 the Carpathian region. *International Journal of Climatology*, 41, E1180-E1201, <https://doi.org/10.1002/joc.6761>,  
798 2021.
- 799 Kling, H., Fuchs, M., & Paulin, M.: Runoff conditions in the upper Danube basin under an ensemble of climate  
800 change scenarios. *Journal of hydrology*, 424, 264-277, <https://doi.org/10.1016/j.jhydrol.2012.01.011>, 2012.
- 801 Kotlarski, S., Keuler, K., Christensen, O. B., Colette, A., Déqué, M., Gobiet, A., Goergen, K., Jacob, D., Lüthi,  
802 D., van Meijgaard, E., Nikulin, G., Schär, C., Teichmann, C., Vautard, R., Warrach-Sagi, K., & Wulfmeyer, V.:  
803 Regional climate modeling on European scales: a joint standard evaluation of the EURO-CORDEX RCM  
804 ensemble. *Geoscientific Model Development*, 7(4), 1297-1333. <https://doi.org/10.5194/gmd-7-1297-2014>, 2014.

- 805 Kotlarski, S., Keuler, K., Christensen, O.B., Colette, A., Déqué, M., Gobiet, A., Goergen, K., Jacob, D., Lüthi,  
806 D., Van Meijgaard, E. and Nikulin, G.: Regional climate modeling on European scales: a joint standard  
807 evaluation of the EURO-CORDEX RCM ensemble. *Geoscientific Model Development*, 7(4), pp.1297-1333.  
808 doi:10.5194/gmd-7-1297-2014, 2014.
- 809 Liemohn, M. W., Shane, A. D., Azari, A. R., Petersen, A. K., Swiger, B. M., & Mukhopadhyay, A.: RMSE is not  
810 enough: Guidelines to robust data-model comparisons for magnetospheric physics. *Journal of Atmospheric and*  
811 *Solar-Terrestrial Physics*, 218, 105624, <https://doi.org/10.1016/j.jastp.2021.105624>, 2021.
- 812 Masanganise, J., Magodora, M., Mapuwei, T., & Basira, K.: An assessment of CMIP5 global climate model  
813 performance using probability density functions and a match metric method. *Science Insights: An International*  
814 *Journal*, 4(1), 1-8, 2014.
- 815 Mirkes, E. M., Allohobi, J., & Gorban, A.: Fractional norms and quasinorms do not help to overcome the curse of  
816 dimensionality. *Entropy*, 22(10), 1105, <https://doi.org/10.48550/arXiv.2004.14230>, 2020.
- 817 Mooney, P. A., Rechid, D., Davin, E. L., Katragkou, E., de Noblet-Ducoudré, N., Breil, M., Cardoso, R. M.,  
818 Daloz, A. S., Hoffmann, P., Lima, D. C. A., Meier, R., Soares, P. M. M., Sofiadis, G., Strada, S., Strandberg, G.,  
819 Toelle, M. H., & Lund, M. T.: Land-atmosphere interactions in sub-polar and alpine climates in the CORDEX  
820 Flagship Pilot Study Land Use and Climate Across Scales (LUCAS) models – Part 2: The role of changing  
821 vegetation, *The Cryosphere*, 16, 1383–1397, <https://doi.org/10.5194/tc-16-1383-2022>, 2022
- 822 Murphy, A. H.: Skill scores based on the mean square error and their relationships to the correlation  
823 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
- 825 Nash, J. E., & Sutcliffe, J. V.: River flow forecasting through conceptual models part I—A discussion of  
826 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.
- 827 Pachepsky, Y. A., Martinez, G., Pan, F., Wagener, T., & Nicholson, T.: Evaluating hydrological model  
828 performance using information theory-based metrics. *Hydrology and Earth System Sciences Discussions*, 1-24,  
829 <https://doi.org/10.5194/hess-2016-46>, 2016.
- 830 Perkins, S. E., A. J. Pitman, N. J. Holbrook, and J. McAneney,: Evaluation of the AR4 Climate Models' Simulated  
831 Daily Maximum Temperature, Minimum Temperature, and Precipitation over Australia Using Probability Density  
832 Functions. *J. Climate*, 20, 4356–4376, <https://doi.org/10.1175/JCLI4253.1>, 2007.
- 833 Pitman, J.: Exchangeable and partially exchangeable random partitions. *Probability theory and related*  
834 *fields*, 102(2), 145-158, <https://doi.org/10.1007/BF01213386>, 1995.
- 835 Randall, D. A., Wood, R. A., Bony, S., Colman, R., Fichefet, T., Fyfe, J., Kattsov, V., Pitman, A., Shukla, J.,  
836 Srinivasan, J., Stouffer, R.J., Sumi, A. & Taylor, K. E.: Climate models and their evaluation. In *Climate Change*  
837 *2007: The physical science basis. Contribution of Working Group I to the Fourth Assessment Report of the IPCC*  
838 *(FAR)* (pp. 589-662). Cambridge University Press, 2007.
- 839 Reich, N. G., Lauer, S. A., Sakrejda, K., Iamsirithaworn, S., Hinjoy, S., Suangtho, P., Suthachana, S., Clapham,  
840 H.E., Salje, H., Cummings, D.A. & Lessler, J.: Challenges in real-time prediction of infectious disease: a case  
841 study of dengue in Thailand. *PLoS neglected tropical diseases*, 10(6), e0004761,  
842 <https://doi.org/10.1371/journal.pntd.0010883>, 2016.
- 843 Richter, J. H., Butchart, N., Kawatani, Y., Bushell, A. C., Holt, L., Serva, F., Anstey, J., Simpson, I.R., Osprey,  
844 S., Hamilton, K., Braesicke, P., Cagnazzo, C., Chen C. C., Garcia, R. R., Gray, L. J., Kerzenmacher, T., Lott, F.,  
845 McLandress, C., Naoe, H., Scinocca, J., Stockdale, T. N., Versick, S., Watanabe, S., Yoshida, K., & Yukimoto,  
846 S.: Response of the quasi-biennial oscillation to a warming climate in global climate models. *Quarterly Journal*  
847 *of the Royal Meteorological Society*, 148(744), 1490-1518, <https://doi.org/10.1002/qj.3749>, 2022.



- 848 Roberts, N. M., & Lean, H. W.: Scale-Selective Verification of Rainfall Accumulations from High-Resolution  
849 Forecasts of Convective Events. *Monthly Weather Review*, 136, 78–97, <https://doi.org/10.1175/2007MWR2123.1>,  
850 2008.
- 851 Rupp, D. E., Abatzoglou, J. T., Hegewisch, K. C., & Mote, P. W.: Evaluation of CMIP5 20th century climate  
852 simulations for the Pacific Northwest USA. *Journal of Geophysical Research: Atmospheres*, 118(19), 10-884,  
853 <https://doi.org/10.1002/jgrd.50843>, 2013.
- 854 Smiatek, G., Kunstmann, H. & Senatore A.: EURO-CORDEX regional climate model analysis for the Greater  
855 Alpine Region: Performance and expected future change, *Journal of Geophysical Research: Atmospheres*, 121,  
856 7710–7728, doi:10.1002/2015JD024727, 2016
- 857 Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram. *Journal of Geophysical  
858 Research: Atmospheres*, 106(D7), 7183-7192, <https://doi.org/10.1029/2000JD900719>, 2001.
- 859 Torma, C. Z.: Detailed validation of EURO-CORDEX and Med-CORDEX regional climate model ensembles  
860 over the Carpathian Region. *Időjárás/Quarterly Journal Of The Hungarian Meteorological Service*, 123(2), 217-  
861 240. DOI:10.28974/idojaras.2019.2.6, 2019.
- 862 Van Noije, T., Bergman, T., Le Sager, P., O'Donnell, D., Makkonen, R., Gonçalves-Ageitos, M., M., Döscher,  
863 R., Fladrich, U., von Hardenberg, J., Keskinen, J.-P., Korhonen, H., Laakso, A., Myriokefalitakis, S., Ollinaho,  
864 P., Pérez García-Pando, C., Reerink, T., Schrödner, R., Wyser, K., & Yang, S.: EC-Earth3-AerChem: a global  
865 climate model with interactive aerosols and atmospheric chemistry participating in CMIP6. *Geoscientific Model  
866 Development*, 14(9), 5637-5668, <https://doi.org/10.5194/gmd-14-5637-2021>, 2021.
- 867 Vautard, R., Kadyrov, N., Iles, C., Boberg, F., Buonomo, E., Bülow, K., et al.: Evaluation of the large EURO-  
868 CORDEX regional climate model ensemble. *Journal of Geophysical Research: Atmospheres*, 126,  
869 e2019JD032344, <https://doi.org/10.1029/2019JD032344>, 2021.
- 870 Wang, Z & Bovik, A. C.: A universal image quality index. *IEEE Signal Processing Letters*, 9, 3, 81-84, doi:  
871 10.1109/97.995823, 2002.
- 872 Watt-Meyer, O., Brenowitz, N. D., Clark, S. K., Henn, B., Kwa, A., McGibbon, J., Perkins, W. A. & Bretherton,  
873 C. S.: Correcting weather and climate models by machine learning nudged historical simulations. *Geophysical  
874 Research Letters*, 48(15), e2021GL092555, <https://doi.org/10.1029/2021GL092555>, 2021.
- 875 Weber, R., Schek, H. J., & Blott, S.: A quantitative analysis and performance study for similarity-search methods  
876 in high-dimensional spaces. In *VLDB*, 98, 194-205, 1998.
- 877 Węglarczyk, S.: The interdependence and applicability of some statistical quality measures for hydrological  
878 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.
- 879 Wehner, M., Lee, J., Risser, M., Ullrich, P., Gleckler, P., & Collins, W. D.: Evaluation of extreme sub-daily  
880 precipitation in high-resolution global climate model simulations. *Philosophical Transactions of the Royal Society  
881 A*, 379(2195), 20190545, <https://doi.org/10.1098/rsta.2019.0545>, 2021.
- 882 Wilcox, R. H.: Adaptive control processes—A guided tour, by Richard Bellman, Princeton University Press,  
883 Princeton, New Jersey, 1961, 255 pp., \$6.50. *Naval Research Logistics Quarterly*, 8(3), 315-316,  
884 <https://www.jstor.org/stable/j.ctt183ph6v>, 1961.
- 885 Willmott, C. J., & Matsuura, K.: Advantages of the mean absolute error (MAE) over the root mean square error  
886 (RMSE) in assessing average model performance. *Climate research*, 30(1), 79-82,  
887 <https://www.jstor.org/stable/24869236>, 2005.
- 888 Wood, R. R., Lehner, F., Pendergrass, A. G., & Schlunegger, S.: Changes in precipitation variability across time  
889 scales in multiple global climate model large ensembles. *Environmental Research Letters*, 16(8), 084022.  
890 <https://doi.org/10.1088/1748-9326/ac10dd>, 2021.

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

894

895

896