

1 Dependence of regionalization methods on the complexity of  
2 hydrological models in multiple climatic regions

3

4 Xue Yang<sup>1</sup>, Jan Magnusson<sup>2</sup>, Shaochun Huang<sup>2</sup>, Stein Beldring<sup>2</sup>, Chong-Yu Xu<sup>1\*</sup>

5

6 <sup>1</sup>Department of Geosciences, University of Oslo, P.O. Box 1047 Blindern, 0316 Oslo, Norway

7 <sup>2</sup>Norwegian Water Resources and Energy Directorate (NVE), P.O. Box. 5091 Majorstua, 0301 Oslo, Norway

8 \* Corresponding to: [c.y.xu@geo.uio.no](mailto:c.y.xu@geo.uio.no)

9 Tel : +0047-22-855825; Fax: +0047-22-854215

10

11 **Abstract** Hydrological models have been widely used to predict runoff in regions with observed  
12 discharge data, and regionalization methods have been extensively discussed for providing runoff  
13 predictions in ungauged basins (PUB), especially during the PUB decade (2003-2012). Great progress  
14 has been achieved in the field of regionalization in previous studies, in which different hydrological  
15 models have been coupled with various regionalization methods. However, different conclusions have  
16 been drawn due to the use of different hydrological models, regionalization methods, and study  
17 regions. In this study, we assessed the performance of the five most widely used regionalization  
18 methods (spatial proximity with parameter averaging option (SP-par), spatial proximity with output  
19 averaging option (SP-out), physical similarity with parameter averaging option (Phy-par), physical  
20 similarity with output averaging option (Phy-out), and regression methods (PCR)) and four daily  
21 rainfall-runoff models (GR4J, WASMOD, HBV and XAJ, with 6, 8, 13, and 19 parameters,  
22 respectively) at the same time. Our aim was to evaluate how the performance of the regionalization  
23 methods depends on (a) the selection of hydrological models, (b) nonstationary climate conditions,  
24 and (c) different climate regions. This investigation used data from 86 independent catchments evenly  
25 distributed throughout Norway, covering three different climate zones (oceanic, continental and polar  
26 tundra) according to the Köppen-Geiger classification. The results showed that (a) the SP-out and Phy-

27 out methods performed better than the SP-par and Phy-par for all the hydrological models, and the  
28 regression method performed worst in most cases; (b) the difference between the parameter averaging  
29 option and the output averaging option is positively related to the number of hydrological model  
30 parameters, i.e. the greater the number of parameters, the larger the difference between the two options;  
31 (c) the XAJ model with the greatest number of parameters produced the best results in most cases, and  
32 models with fewer parameters tend to produce similar performance for the different regionalization  
33 methods; (d) models with more parameters displayed larger declines in performance than those with  
34 fewer parameters for nonstationary conditions; and (e) clear differences in the performance of the  
35 regionalization methods exist among the three climate regions. This study provides insight into the  
36 relationship between the complexity of hydrological models and regionalization methods in cold and  
37 seasonally snow-covered regions.

38 **Keywords:** Regionalization methods; hydrological models; climate zones; ungauged catchments

1 Dependence of regionalization methods on [the](#) complexity of  
2 hydrological models in multiple [climatic](#) regions

3  
4 **Abstract** Hydrological models have been widely used to predict runoff in regions with observed  
5 discharge data, and regionalization methods have been extensively discussed for providing runoff  
6 predictions in ungauged basins (PUB), especially during the PUB decade (2003-2012). Great progress  
7 has been achieved in the field of regionalization in previous studies, in which different hydrological  
8 models have been coupled with various regionalization methods. However, different conclusions have  
9 been drawn due to the use of different hydrological models, regionalization methods, and study  
10 regions. In this study, we assessed the performance of the five most widely used regionalization  
11 methods (spatial proximity with parameter averaging option (SP-par), spatial proximity with output  
12 averaging option (SP-out), physical similarity with parameter averaging option (Phy-par), physical  
13 similarity with output averaging option (Phy-out), and regression methods (PCR)) and four daily  
14 rainfall-runoff models (GR4J, WASMOD, HBV and XAJ, with 6, 8, 13, and 19 parameters,  
15 respectively) at the same time. Our aim was to evaluate how the performance of the regionalization  
16 methods depends on (a) the selection of hydrological models, (b) nonstationary climate conditions,  
17 and (c) different climate regions. This investigation used data from 86 independent catchments evenly  
18 distributed throughout Norway, covering three different climate zones (oceanic, continental and polar  
19 tundra) according to the Köppen-Geiger classification. The results showed that (a) [the](#) SP-out and Phy-  
20 out methods performed better than the SP-par and Phy-par for all the hydrological models, and the  
21 regression method performed worst in most cases; (b) the difference between the parameter averaging  
22 option and the output averaging option is positively related to the number of hydrological model  
23 parameters, i.e. [the greater the](#) number of parameters, the larger the difference between the two options;  
24 (c) the XAJ model with [the greatest](#) number of parameters produced the best results in most cases, and  
25 models with fewer parameters tend to produce similar performance for the different regionalization  
26 methods; (d) models with more parameters displayed larger declines in performance than those with  
27 fewer parameters for nonstationary conditions; and (e) clear differences in the performance of the

28 regionalization methods exist among the three climate regions. This study provides insight into the  
29 relationship between the complexity of hydrological models and regionalization methods in cold and  
30 seasonally snow-covered regions.

31

32 **Keywords:** Regionalization methods; hydrological models; climate zones; ungauged catchments

33

## 34 **1. Introduction**

35

36 Runoff prediction plays a significant and essential role in water resources management, the assessment  
37 of the impact of environmental change (e.g., climate and land use), and hydrological design (e.g.,  
38 Blöschl and Montanari, 2010; Parajka et al., 2013). During the last several decades, hydrological  
39 models have become the most popular and common solution for runoff predictions. However, the  
40 models have free parameters to be calibrated by using the observed discharge data before predicting  
41 the runoff hydrographs, which are not available in many catchments of interest (e.g., He et al., 2011;  
42 Parajka et al., 2013). This fact made the topic ‘predictions in basins without observed discharge data  
43 (ungauged basins)’ attractive and challenging for hydrologists (e.g., Parajka et al., 2007; Sivapalan et  
44 al., 2003; Xu, 2003). As a result, the International Association of Hydrological Sciences (IAHS)  
45 established a “Decade on Predictions in Ungauged Basins (PUB): 2003–2012”, and great progress has  
46 been achieved during this period (Hrachowitz et al., 2013).

47

48 Regionalization is defined as the method for predicting runoff in ungauged basins by transferring  
49 information from gauged (donor) to ungauged (target) catchments (e.g., Rojas- Serna et al., 2016;  
50 Razavi and Coulibaly, 2013). In general, regionalization methods are classified into three categories:  
51 (a) spatial proximity methods assume that geographically close catchments have similar hydrological  
52 behaviors (e.g., Egbuniwe and Todd, 1976; Vandewiele et al., 1991); (b) physical similarity  
53 methods assume that catchments with similar physical characteristics have the same hydrological

54 response (e.g., Burn and Boorman, 1993; McIntyre et al., 2005), thus, the parameter values are  
55 transferred to ungauged basins from either geographically close or physically similar gauged basins;  
56 and (c) the regression method, which is one of the most popular and oldest regionalization approaches  
57 (Oudin et al., 2008), links model parameters to physical and climatic catchment characteristics by  
58 regression functions and assumes that the relationship is transferable from gauged to ungauged basins  
59 (e.g., Magette et al., 1976; Young, 2006).

60

61 Many studies have applied and compared regionalization methods for various regions in combination  
62 with a wide range of hydrological models. However, in many cases, the conclusion about which  
63 method performed best differs largely among the studies. For example, Merz and Blöschl (2004)  
64 concluded that the spatial proximity method performed better than the regression method for  
65 catchments in Austria using the HBV model. On the other hand, Young (2006) found that the  
66 regression method gave better results than the spatial proximity method in the UK. Bao et al. (2012)  
67 concluded that the physical similarity method was best by using the Akaike information criterion (AIC)  
68 on 55 catchments in China. Different models were applied for different regions in these studies, and  
69 therefore many hydrologists claim that the performance of regionalization methods depends on the  
70 study area and the choice of hydrological model (e.g., Parajka et al., 2013; Reichl et al., 2009; Salinas  
71 et al., 2013; Samuel et al., 2011; Viglione et al., 2013). Most of the above-mentioned studies only used  
72 one hydrological model in a specific region, and conclusions cannot be drawn on how the model  
73 selection or study region affects the performance of the regionalization methods.

74

75 Few studies have assessed the performance of regionalization methods using multiple models. Li et al.  
76 (2017) used SIMHYD (10 model parameters) and XAJ (12 model parameters) in Australia and found  
77 consistent regionalization results for both models. The same conclusion was drawn by Li et al. (2014),  
78 where GR4J (7 model parameters) and SIMHYD (12 model parameters) were applied in the southeast  
79 Tibetan Plateau. Furthermore, Petheram et al. (2012) conducted a comparison by using five rainfall-  
80 runoff models and concluded that the difference between hydrological models was negligible for  
81 runoff prediction in ungauged basins. This conclusion was consistent with two other studies (Chiew,

82 2010; Viney et al., 2009b), which also included five hydrological models. However, none of these  
83 studies included a regression approach, which provided very different results when used with either  
84 the GR4J (4 model parameters) or TOPMO (6 model parameters) model in the study of Oudin et al.  
85 (2008), who tested three kinds of regionalization methods using two hydrological models for 913  
86 catchments in France. Either the number of regionalization methods or the number of models used in  
87 previous studies is still too small to draw a general conclusion. In addition, all these evaluations have  
88 been performed for relatively warm climate regions, where the snow process is of limited importance.  
89 Thus, a more comprehensive study is needed to investigate how regionalization performance differs  
90 with multiple hydrological models of different complexity for runoff prediction in ungauged basins,  
91 especially for cold and seasonally snow-covered regions.

92

93 Furthermore, climate is changing (IPCC, 2014), resulting in nonstationary relationships between  
94 rainfall and runoff (Zhang et al., 2011a), which makes the reliability of applying the conclusions made  
95 in a historical period into future application questionable. Thus, for future runoff prediction in  
96 ungauged basins, it is essential to investigate the transferability of the regionalization methods under  
97 changing climatic conditions (e.g., Broderick et al., 2016; Yang et al., 2019). Finally, regionalization  
98 performances also vary between regions, according to Parajka et al. (2013), who statistically  
99 summarized this conclusion from 34 regionalization studies. However, it cannot explicitly present the  
100 performance difference between regions for specifically selected regionalization methods because  
101 different hydrological models and regionalization methods were applied in the studies cited and  
102 summarized by Parajka et al. (2013).

103

104 In this study, we perform a comprehensive evaluation of the performance of five widely used  
105 regionalization methods (see section 3.2) combined with four frequently used hydrological models  
106 (GR4J–6 parameters, WASMOD–8 parameters, HBV–13 parameters and XAJ–19 parameters) in  
107 regions with highly contrasting physiographic and climatic settings. The evaluation is based on 86  
108 catchments in Norway, belonging to three different climatic regions according to the Köppen-Geiger  
109 classification (Kottek et al. 2006) and under different climate conditions. This is the first study that

110 specifically addresses how the performance of the regionalization methods (a) depends on the  
111 selection of hydrological models, (b) changes in different climate conditions, i.e., when air  
112 temperature increases, and (c) varies between different climate regions as defined by the Köppen-  
113 Geiger classification.

114

## 115 **2. Study area and data**

### 116 *2.1 Study area*

117 Our study catchments are located in Norway, which is situated in northern Europe in the western and  
118 northern part of the Scandinavian Peninsula. Norway has a long and rugged coastline, elevation  
119 [spanning](#) from sea level to 2469 m.a.s.l., and latitudes [ranging](#) from 58° to 71°N. This results in highly  
120 variable hydroclimatological conditions across the study domain (Vormoor et al., 2016; Yang et al.,  
121 2018, 2019). In this study, we used data from 86 nonoverlapping catchments distributed evenly  
122 throughout our study domain (Figure 1). These stations have continuous meteorological data and  
123 discharge data records with less than 40% missing values during the periods from 1980 to 1989 as  
124 well as 2006 to 2015. These two periods are used in this study. The left panel map in [Figure. 1](#) also  
125 displays the Köppen-Geiger climate classification, which is based on data from 1976 to 2000 (Kottek  
126 et al. 2006; Peel et al., 2007; Beck et al., 2018). Note that the original classification divided Norway  
127 into five different climate groups. However, in two of these groups, less than 10 catchments were  
128 located. We therefore merged some of the groups, resulting in the following three regions: (a) oceanic  
129 climate containing 19 catchments, (b) continental climate containing 52 catchments and (c) polar  
130 tundra climate containing 15 catchments.

131

132

Fig 1. Insert here

133

134 **2.2 Data**

135 For the hydrological simulations, we used daily precipitation and temperature data acquired from the  
136 gridded seNorge dataset with a resolution of 1 km produced by the Norwegian Meteorological  
137 Institute (Tveito et al., 2005; Mohr, 2009; Jansson et al., 2007). Daily discharge data were obtained  
138 from the hydrometric observation network of the Norwegian Water Resources and Energy Directorate  
139 (NVE). To test the performance of the regionalization methods under varying climate conditions, we  
140 analyzed the precipitation and temperature records for the period from 1980 to 2015 (Figure 2). For  
141 precipitation, there is no clear trend, whereas temperature increases throughout the study period. For  
142 model calibration and verification, we selected ten years at the start (1980 to 1989) and the end (2006  
143 to 2015) of the whole period since these two periods show the largest difference in air temperature.  
144 For the first period, the average precipitation is 1932 mm/year, and the air temperature is 1.2°C. For  
145 the second period, the average precipitation is 2027 mm/year, and the air temperature is 2.6°C. The  
146 right panels in Figure 2 show the average monthly precipitation, temperature and Pardé coefficient  
147 (ratio between the average monthly discharge and the mean annual runoff) for the catchments in each  
148 climate group. The oceanic climate group is characterized by higher precipitation during autumn and  
149 winter and higher air temperature than that of the two remaining groups. The watersheds in the oceanic  
150 climate group also show two peaks in runoff (compare the Pardé coefficient between the groups)  
151 resulting from spring snowmelt and strong rainfall during autumn. The continental climate group  
152 displays low seasonality for precipitation but high seasonal variations in temperature, resulting in one  
153 peak runoff caused by snowmelt. The climate characteristics for the polar tundra climate group are  
154 similar to those of the continental group, but with lower temperature, and the snowmelt-induced peak  
155 in runoff occurs later.

156

157

158

Fig 2. Insert here.

159 Table 1 shows the average annual and seasonal precipitation, temperature and runoff for the three  
160 climate classes. Precipitation in the oceanic climate group is substantially larger than that in the other  
161 two groups, which show rather similar precipitation amounts. For temperature, the oceanic climate



162 group shows the highest values, whereas the coldest temperatures are recorded in the polar tundra  
163 climate group. In particular, for the oceanic group, precipitation increases from the calibration to  
164 verification period for the winter season, but for the summer season, the difference is small between  
165 the two periods. For temperature, the increase from the calibration to verification period is smallest in  
166 the oceanic region compared to the other regions. The seasonal characteristics in runoff are similar to  
167 those of precipitation. Note that summer runoff decreases from the calibration to the verification  
168 period for all groups.

169

170 Table 1 Insert here.

171

172 Since there is no potential evapotranspiration ( $E_p$ ) data available in our study area, which are needed  
173 as the input data for the hydrological models, we applied the Hargreaves equation (Hargreaves, 1975)  
174 to calculate  $E_p$  (mm/day), which is recommended by Shuttleworth (1993) and Xu et al. (2002):

$$175 \quad E_p = 0.0023 R_a (TC + 17.8)\sqrt{TR} \quad (1)$$

176 where  $R_a$  is the extraterrestrial radiation for the location in mm/day evaporation equivalent (Allen et  
177 al., 1998),  $TC$  is the temperature ( $^{\circ}C$ ), and  $TR$  is the daily temperature range ( $^{\circ}C$ ).

178

179 A set of catchment descriptors is needed for two of the regionalization methods, namely, the physical  
180 similarity and regression methods (see Table 2). These catchment descriptors were used in Yang et al.  
181 (2018, 2019). Similar catchment descriptors have been used in several studies for evaluating  
182 regionalization methods (e.g., He et al., 2011; McIntyre et al., 2005; Merz and Blöschl, 2004).

183

184 Table 2 Insert here.

185

## 186 **3. Methods**

### 187 *3.1 Hydrological models*

188 Four widely used conceptual rainfall-runoff models running at a daily time step were selected for the  
189 analysis in this study, and a snow module was included in the models since runoff in many of the  
190 catchments is strongly affected by the accumulation and melting of snow. The number of model  
191 parameters varies from 6 to 17 between the models after adding the snow routine. Figure 3 shows the  
192 model structures, and a description of the parameters is available in Table 3.

193

194 GR4J (Génie Rural à 4 paramètres Journalier) is a model based on unit hydrograph principles with  
195 four free parameters (Perrin et al., 2003). It has been widely used in regionalization studies worldwide,  
196 such as in France (Oudin et al., 2008), China (Li et al., 2014) and Australia (Zhang et al., 2014, 2016).  
197 We coupled the GR4J model with a degree-day type snow module called CemaNeige that was  
198 developed by Valéry (2010). This snow module allows us to estimate snowmelt and simulate  
199 snowpack evolution using 2 additional parameters, and the coupling of GR4J and CemaNeige has  
200 been tested in other studies (e.g., Coron et al., 2014; Hublart et al., 2015).

201

202 WASMOD (The Water And Snow balance modelling system) is a model with simple structure and has  
203 been validated in many different climate regions (e.g., Xu and Singh, 2002; Li et al., 2013, 2015;  
204 Widén-Nilsson et al., 2007; Xu and Halldin, 1997). For regionalization studies, it has been applied in  
205 Sweden (Xu, 2003), Denmark (Muller-Wohlfeil et al., 2003) and Norway (Yang et al., 2018; 2019).  
206 The version of WASMOD used in this study has eight free parameters.

207

208 HBV (Hydrologiska Byråns Vattenbalansavdelning) is a popular model used for runoff simulation in  
209 both gauged and ungauged basins. For regionalization studies, it has been applied in different climate  
210 regions, such as Austria (e.g., Merz and Blöschl, 2004; Parajka et al., 2005), Sweden (Seibert and  
211 Beven, 2009), China (Jin et al., 2009), Canada (Samuel et al., 2011) and the US (Pool et al., 2017). In  
212 our study, we followed the structure and formulas in the HBV-light version (Seibert and Vis, 2012),

213 which includes a snow routine, soil moisture routine, response function and routing routine. In total,  
214 this model has 13 calibration parameters.

215

216 The XAJ (Xin An Jiang) model was developed for humid regions in China by Zhao et al. (1980, 1992)  
217 and has since become a widely used model in flood forecasting, water resources assessment, and  
218 climate change assessments. The original model consists of modules for computing evapotranspiration,  
219 runoff production, runoff separation, and flow routing. It has also been applied in many  
220 regionalization studies (e.g., Zhang and Chiew, 2009; Li et al., 2009, 2017). We implemented the  
221 structure shown in Lin et al. (2014) without the Muskingum routing module because our catchments  
222 are rather small in size with steep slopes, and therefore, river flow routing is not an important process  
223 (Li et al., 2014). However, there is no snow module in XAJ, and therefore, we coupled it with the  
224 CemaNeige snow module (see description of the GR4J model above). This model system contains 17  
225 parameters in total.

226

227 Fig 3. Insert here

228

229 Table 3 Insert here.

230

### 231 **3.2 Regionalization methods**

232 Spatial proximity, physical similarity and regression methods are commonly used in regionalization  
233 studies (e.g., Oudin et al., 2008; Petheram et al., 2012; Hrachowitz et al., 2013). For spatial proximity  
234 and physical similarity methods, which are classified as distance-based regionalization methods  
235 according to He et al. (2011), the model parameter values in ungauged catchments are transferred from  
236 gauged donor catchments. For the regression method, the model parameter values in ungauged  
237 catchments are determined by regression functions established using data from gauged basins. The  
238 regression method in this study is principal component regression (PCR), which couples principal  
239 component analysis (PCA) with the multiple linear regression method. Using PCA, a set of  
240 observations of possibly correlated catchment descriptors is converted into a set of linearly  
241 uncorrelated variables called principal components. Then, the relationships among model parameters

242 and selected catchment descriptors are established using multiple linear regression. Finally, the  
243 functions are used for estimating model parameters in the ungauged catchments. Table 4 describes the  
244 equations and assumptions for the regionalization methods applied in this study.

245

246 Table 4 Insert here.

247

248 For distance-based regionalization methods, i.e., spatial proximity and physical similarity, two  
249 approaches are often used for transferring the model parameters from the gauged donor to the  
250 ungauged target catchments (e.g., McIntyre et al., 2005; Oudin et al., 2008): (a) for the so-called  
251 parameter averaging option, the model parameters from the donor catchments are first averaged and  
252 then used to run the model for the target catchment, and (b) for the so-called output averaging option,  
253 the model is first run using the parameter sets from the donor catchments (i.e., basins with runoff  
254 where model calibration is possible) on the target catchment and the outputs from the model are then  
255 averaged. As a result, there are five regionalization approaches used in this study, as shown in Table 5.  
256 For a more detailed description and similarity index introduction, please see Yang et al. (2018, 2019).

257

258 Table 5 Insert here.

259

### 260 **3.3 Performance evaluation**

#### 261 3.3.1 Model calibration and verification

262 In this study, we applied a widely used objective function proposed by Viney et al. (2009a) when  
263 calibrating the models. This objective function is a weighted combination of the Nash and Sutcliffe  
264 efficiency (Nash and Sutcliffe, 1970) and a logarithmic penalty function based on the bias as follows:

$$265 \quad F = NSE - 5 * |\ln(1 + bias)|^{2.5} \quad (2)$$

266 where:

$$267 \quad NSE = 1 - \frac{\sum(Q_{sim} - Q_{obs})^2}{\sum(Q_{obs} - \overline{Q_{obs}})^2} \quad (3)$$

$$268 \quad bias = \frac{\overline{Q_{sim}} - \overline{Q_{obs}}}{\overline{Q_{obs}}} \quad (4)$$

269  $Q_{\text{obs}}$  represents the observed runoff, and  $Q_{\text{sim}}$  represents the simulated runoff. F values can vary from  
270  $-\infty$  to the optimal value of 1. This objective function can come close to maximizing Nash and Sutcliffe  
271 efficiency (NSE) and minimizing the bias at the same time (Vaze et al., 2010). For the calibration  
272 process, we used a standard gradient-based automatic optimization method (Lagarias et al., 1998)  
273 implemented in the MATLAB software package (“fmincon” function; MATLAB R2016b, The  
274 MathWorks, Inc., Natick, Massachusetts, United States).

275

276 The split-sample test is commonly used for model verification, aiming to show the model validity in  
277 different climate conditions (e.g., Coron et al., 2012; Xu, 1999; Klemeš, 1986). In the current study,  
278 we evaluate the model performance for 1980-1989 and 2006-2015, and the temperature and  
279 precipitation in the latter period are approximately 1.4°C and 5% higher than that in the first period.

280

### 281 3.3.2 Evaluation of regionalization methods

282

283 We performed three different evaluations of the regionalization methods. In the first evaluation, the  
284 performance of the regionalization methods was tested for all models using data from the calibration  
285 period, aiming to show the differences among the models. In this step, we applied a leave-one-out  
286 cross verification method as in many other studies (e.g., Yang et al., 2018; McIntyre et al., 2005). In  
287 the second analysis, we repeated the same evaluation but for the warmer and wetter verification period.  
288 This analysis thus tests the transferability of both the regionalization methods and hydrological models  
289 under climate change conditions (e.g., Broderick et al., 2016; Li et al., 2012). In the final evaluation,  
290 we summarize and discuss the performance of the regionalization methods for the three different  
291 climatic regions (see section 2.1). Since the climate is changing to be warmer in the future (IPCC,  
292 2014), the following regionalization performance for different climate conditions is investigated from  
293 1980-1989 (calibration) to 2006-2015 (verification).

294

### 295 3.3.3 Evaluation criteria

296 To investigate the performance from different aspects, we applied four different criteria in this study.  
297 The calibration function F (Equation 2) is the first selection since it considers both the goodness of fit  
298 and the water balance aspects between the simulated and observed runoff. NSE (Equation 3) is the most  
299 commonly used criterion in hydrology to measure the fit of the hydrographs between the observed and  
300 simulated runoff, which is relatively sensitive to high flow (e.g., Oudin et al., 2008; Pushpalatha et al.,  
301 2012; Zhang and Chiew, 2009). Similarly, we included another criterion, NSElog, which is based on  
302 the same formulation as NSE but computed on logarithmic transformed flows and with more emphasis  
303 on low flow (e.g., Oudin et al., 2008; Pushpalatha et al., 2012). Finally, the percentage of bias (Pbias)  
304 (Equation 4) is applied to measure the average tendency of the simulation to be larger or smaller than  
305 the observed counterparts.

306

307 The range for F, NSE and NSElog is  $(-\infty, 1)$ , where 1 means the simulated runoff perfectly fits the  
308 observed runoff and less than 0 suggests that the model is no better than the observed mean value. For  
309 Pbias, it varies between  $(-\infty, +\infty)$  with the optimal value equal to 0 and worse performance for water  
310 balance simulation if the absolute Pbias is larger.

## 311 **4. Results**

### 312 *4.1 Hydrological model performance in cross verification*

313 Before evaluating both the hydrological models and the regionalization methods, we first assessed the  
314 performance of the models by a split-sample test. Figure 4 presents the cumulative density function  
315 (CDF) curves for all hydrological models over 86 catchments, measured by F value during 1980-1989  
316 and 2006-2015.

317

318 For the first calibration period 1980 – 1989 (the left panel in Fig. 4), the CDF curves from all the  
319 hydrological models stay close, and XAJ appears to be slightly better. The average F value is  
320 approximately 0.75 for XAJ, 0.73 for WASMOD, 0.72 for HBV and 0.69 for GR4J. In the verification

321 period 2006 - 2015, the models perform differently, meaning the temporal transferability varies  
322 between the hydrological models. However, the best performance is still produced by XAJ, whose  
323 mean F value is approximately 0.68, followed by WASMOD (0.64). The HBV model shows the worst  
324 performance, with a mean F value of approximately 0.61 and the highest degradation of performance  
325 between the calibration and verification periods.

326  
327 The results in the right panel (calibration in 2006-2015 and verification in 1980-1989) shows very  
328 similar characteristics to those in the left panel. XAJ produced the best performance for both the  
329 calibration and the verification periods. Following the rating classification from Moriasi et al. (2007),  
330 who labeled the performance as ‘good’ if NSE is larger than 0.65 and |Pbias| is less than 15%, the F  
331 values larger than 0.61 are considered “good” model performance. Considering the average aspect, all  
332 mean F values for our split-sample test are higher than 0.61. Thus, all hydrological models applied in  
333 the current study are classified as ‘good’ performing models for runoff simulation for both calibration  
334 and verification periods.

335  
336  
337

Fig 4. Insert here.

338 Table 6 gives the average model performance corresponding to the split-sample test by using other  
339 assessment criteria. First, regarding the water balance aspect, all models yield similarly ‘good’  
340 performance for both subperiods with |Pbias| values smaller than 5%. Second, the model performance  
341 measured by NSE shows consistent findings with the results from the F value, i.e., (a) the models  
342 show similar performance in the calibration period but perform differently in the verification period;  
343 (b) XAJ is considered the best-performing model for both the calibration and the verification cases;  
344 and (c) HBV shows the largest decline in performance from the calibration to the verification period.  
345 This similarity between the results from the F value and NSE can be explained by the small Pbias for  
346 all the simulation results. Finally, according to the results of NSElog, which is more sensitive to low  
347 flow, the simple models (GR4J and WASMOD) display higher values in the calibration period, while  
348 WASMOD and XAJ show better performance in the verification period. Considering the performance

349 loss from calibration to verification, relatively larger degradation appears for the NSElog than for the  
350 NSE and Pbias, especially for the GR4J model.

351

352

Table 6 Insert here.

353

## 354 ***4.2 Evaluation of regionalization methods***

355 4.2.1 Influence of the number of donor catchments on performance under stationary  
356 conditions

357 Figure 5 shows that the output averaging option gives better average performance than the parameter  
358 averaging option in both spatial proximity and physical similarity methods and for all the models,  
359 except for the case of one donor catchment, where both options provided the identical results as  
360 expected. When considering the number of donor catchments, the largest increase in performance  
361 typically occurs when changing from using one donor catchment to using two donor catchments, with  
362 the parameter option for XAJ as the only exception. This is in line with earlier studies that the number  
363 of donor catchments typically affects the performance of distance-based regionalization methods (e.g.,  
364 Oudin et al., 2008; Yang et al., 2018). However, the number of donor catchments providing the best  
365 performance differs among the hydrological models and regionalization methods. For instance, for  
366 XAJ, two donor catchments give the best results for SP-out, whereas 8 donor catchments are needed  
367 for HBV to achieve the optimal performance. Finally, the difference in performance between the  
368 output and parameter averaging options increases with the number of model parameters. For example,  
369 the difference in the average F value between the two options for the GR4J model was approximately  
370 0.025 and increased to 0.075 for XAJ. Thus, when using a model with many parameters, it is more  
371 important to use the output averaging option to achieve optimal performance for runoff simulations in  
372 ungauged basins.

373

374

Fig 5. Insert here.

375



376 The physical similarity methods require fewer donor catchments to achieve optimal performance for  
377 runoff simulations in ungauged basins compared to [that for](#) the spatial proximity methods (Table 7).  
378 On average, the best performance by the physical similarity methods was produced by 3 donor  
379 catchments, whereas the corresponding number for the spatial proximity methods was 8. It is also  
380 [noteworthy](#) that the parameter averaging option [requires fewer](#) donor catchments than the output  
381 averaging option for both [the](#) physical similarity and [the](#) spatial proximity methods. Therefore, for  
382 practical applications, it is highly recommended to analyze the relationship between the  
383 regionalization performance and [the](#) number of donor catchments to choose the best configuration to  
384 obtain [the](#) optimal results for each case.

385

386

Table 7. Insert here.

387

#### 388 4.2.2 Regionalization performance assessment for all catchments

389 As discussed in section 2.2 (Figure 2 and Table 1), the climate conditions, especially air temperature,  
390 differed between 1980-1989 and 2006-2015. This section presents the influence [of](#) climate conditions  
391 on regionalization performance when the models are calibrated in 1980-1989. The evaluation [results](#)  
392 presented here applied the optimized number of donor catchments for each method and model, as  
393 shown in Table 7.

394

#### 395 *Comparison of regionalization performance between hydrological models*

396

397 Figure 6 shows the distribution of F values as split violin plots for the five regionalization methods  
398 and four hydrological models for both the calibration and verification periods. Foremost, for all [the](#)  
399 hydrological models, the regionalization methods applying the output averaging option (SP-out and  
400 Phy-out) showed better performance than the parameter averaging option (SP-par and Phy-par), and  
401 [the](#) regression method is [the](#) worst (compare black dots with circles). This ranking applies for [both](#) the  
402 calibration and [the](#) verification periods, where the methods with output averaging [options](#) presented  
403 more negative skewed distributions and higher mode values than [those of](#) the other [methods](#). On the

404 other hand, for both periods, the difference in the average performance between the regionalization  
405 methods is smaller for GR4J than for the other models. This difference seems to increase with the  
406 number of model parameters and is thus largest for XAJ. For instance, in the calibration period, the  
407 range in the average F values between the regionalization methods equals 0.04 for GR4J and 0.09 for  
408 XAJ. Finally, from the calibration to verification period, performances decreased for all the  
409 hydrological models and regionalization methods but to various extents. Measured by the decrease in  
410 the overall mean F values from the calibration (solid line) to verification (dashed line) period, HBV  
411 and XAJ displayed larger declines in performance than those of GR4J and WASMOD.

412

413

Fig 6. Insert here.

414

415 Figure 7 compares the regionalization performance in terms of the average values of Pbias, NSE and  
416 NSElog for all catchments using four hydrological models in the calibration and verification periods.  
417 Appendix A presents the violin plot for all the evaluation criteria over all the tested catchments.

418

419 Regarding the water balance simulation, all average values of Pbias vary within (-10%, 10%). The  
420 smallest water balance error for regionalized runoff simulation varies with the hydrological models  
421 and regionalization methods. In general, SP-out and Phy-out tend to yield smaller errors for water  
422 balance simulation than those of the other methods.

423

424 The NSE results give similar findings as the F value. First, SP-out and Phy-out methods perform best  
425 for all the hydrological models, with all average NSE values larger than 0.6, and PCR performs worst.  
426 Second, the difference in NSE between the regionalization methods increases with the growing  
427 number of parameters for the hydrological models. For example, the regionalization performance in  
428 the calibration period ranges within (0.57, 0.61) for GR4J and (0.57, 0.67) for XAJ. Third, relatively  
429 larger degradation of the average regionalization performance is found using the HBV and XAJ  
430 models from the calibration to the verification period.

431

432 For the low-flow evaluation, the regionalization methods with the output average option (SP-out and  
433 Phy-out) substantially outperform the other methods, and the performance differences between the  
434 regionalization methods are more distinct for HBV and XAJ. Furthermore, the average performance of  
435 the regionalization methods is highly influenced by the hydrological models. In this study, WASMOD  
436 and HBV produced the highest and lowest average NSElog values for the regionalization methods,  
437 respectively. Compared with the results from the NSE and F values, the evaluation by NSElog  
438 presents a more recognizable performance difference between the regionalization methods and  
439 hydrological models, as well as the difference between the two subperiods.

440

441 Fig 7. Insert here.

442

443 *Comparison of performance between regionalization methods*

444

445 Figure 8 compares the performance difference in terms of NSE and NSElog between the hydrological  
446 models for each regionalization method during the calibration and verification periods. We omit the  
447 results of the F value and Pbias in the following analysis due to high similarity between the results  
448 from the F value and NSE (see Figure 6 and Appendix A) and small average |Pbias| values (see Figure  
449 7).

450

451 According to the average NSE values, XAJ is considered the best hydrological model for all the  
452 distance-based regionalization methods and the second best model for PCR. GR4J shows the best  
453 results for PCR, but the difference in performance between the models (the gray bars for PCR) is  
454 smallest among the regionalization methods, indicating that the hydrological models have relatively  
455 smaller influence on the regression method than on the distance-based methods. However, this  
456 difference is enhanced from the calibration to the verification period, indicating a larger influence of  
457 the hydrological model on future runoff predictions. According to NSElog, WASMOD shows the best  
458 performance for all the regionalization methods and for both periods. In general, a larger difference  
459 between the hydrological models appears for low flows (indicated by NSElog) than for high flows  
460 (indicated by NSE).

461

462

Fig 8. Insert here.

463

#### 464 4.2.3 Assessment of regionalization performance for different climatic regions

465 The three climate regions shown in Figure 1 display very different runoff regimes, particularly

466 between the oceanic and the two remaining groups (Figure 2). For illustration purposes, the

467 dependence of the performance of the regionalization methods on the geographical regions as

468 measured by NSE is shown in Figure 9. It is seen that the oceanic region presented generally better

469 regionalization performance than that of the other two regions, whose performance variation was

470 smaller as well (only four performance classes shown on the figure). Then, some common

471 characteristics are presented in all the regions. First, when considering the regionalization methods, the

472 output averaging option tended to give higher performance than all the other methods. When focusing

473 on the hydrological models, XAJ showed the best performance in most cases for both the calibration

474 and verification periods. Otherwise, none of the remaining models consistently showed better results

475 than the other models for all climate regions and regionalization methods. Finally, GR4J produced the

476 lowest variation in performance within the climate regions between the regionalization methods in

477 almost all cases. From the calibration to verification period, the highest ranking for XAJ with SP-out

478 and Phy-out methods did not change.

479

480

Fig 9. Insert here.

481

## 482 **5. Discussion**

### 483 *5.1 Hydrological model performance*

484 According to the performance classification presented by Moriasi et al. (2007), the split-sample test

485 result in our study indicated that all the hydrological models were able to provide ‘good’ simulations

486 of runoff for both the calibration and the verification periods. Especially for the water balance

487 simulation, the mean values of |Pbias| for all the studied models are smaller than 5%.

488

489 According to the evaluations in the calibration period based on the F value and NSE in our study area,  
490 XAJ is the best-performing model, and the performance tends to decrease with a decrease in the  
491 number of parameters for the hydrological models. This finding is in line with the statement that  
492 increasing the number of model parameters can lead to better performance during the calibration  
493 period (e.g., Perrin et al., 2001; Petheram et al., 2012; Parajka et al., 2013). However, the result in  
494 terms of low flow simulation (evaluations by NSElog) did not support that statement. For example,  
495 WASMOD outperformed XAJ and HBV for both subperiods. Therefore, further study is needed to  
496 assess the relationship between hydrological model complexity and performance in terms of low flow.  
497 Furthermore, for the verification results, the performances among the models varied substantially. The  
498 degradation of performance is quite similar between the hydrological models evaluating by the F value  
499 and NSE, but distinct differences are shown in the NSElog results. It reminds us that specific criteria  
500 are needed for evaluation of hydrological models when the emphasis stands on low flow or draughts.  
501 Regarding the model performance change from the calibration to the verification period, the model  
502 performance of the XAJ model did not vary substantially. This is incompatible with earlier findings,  
503 which suggest that a complex model tends to have less stable performance than simple models in the  
504 verification period (e.g., Perrin et al., 2001; Holländer, 2009). This phenomenon might relate to the  
505 model structure; for instance, the runoff concentration in the XAJ model includes surface runoff,  
506 interflow runoff and groundwater runoff with three parameters that may better represent the processes  
507 in our study catchments.

508

## 509 ***5.2 Evaluation of regionalization methods***

### 510 **5.2.1 Influence of the number of donor catchments on performance**

511

512 To test the influence of the number of donor catchments on model performance, we examined the  
513 relationship between regionalization performance and the number of donor catchments for all the  
514 models with distance-based methods. The results indicate that using one donor catchment, which

515 might be either the spatially nearest or physically most similar watershed, gives worse results than  
516 using a set of donor catchments. This conclusion is supported by all the tested models in our study,  
517 which is in line with previous findings (e.g., Arsenault and Brissette, 2014; Oudin et al., 2008).  
518 Multiple donor catchments typically provide more information than single donor catchments, which  
519 may explain the behavior described above (e.g., Viney et al., 2009b). However, the output averaging  
520 option might tend to smooth the flow variability as the number of donor catchments increases. This is  
521 especially the case if the donors give models with different time lags between rainfall and peak flow.  
522 Therefore, the smoothing effect and trade-off between the benefits of gains in performance with "more  
523 information" and loss of performance due to this possible smoothing is worth further investigation in  
524 future studies. Our results additionally confirmed that the output averaging option provided better  
525 performance than the parameter averaging option in all the model and method combinations (e.g.,  
526 Oudin et al., 2008, Bao et al., 2012; Yang et al., 2018). Since we applied hydrological models with  
527 different complexities and number of parameters, a promising and new finding is presented in this  
528 study: the difference in performance between the parameter averaging and output averaging options  
529 increases with the number of model parameters (see Figure 5). First, this result can be explained by the  
530 'nonlinear independence' influence between model parameters; thus, transferring the linearly  
531 interpolated individual model parameter value (the parameter averaging option) will lead to  
532 unreasonable model parameters and results (Bárdossy, 2007). Second, hydrological models with more  
533 parameters tend to increase the interaction between their parameters (e.g., Perrin et al., 2003; Poissant  
534 et al., 2017). Hence, we should consider the model parameters as a whole set rather than individual  
535 values for regionalization research as suggested by Bárdossy (2007) and Oudin et al. (2008).

536

537 Some previous studies used one donor catchment for regionalization evaluation according to spatial or  
538 physical similarity and concluded that the difference in performance between hydrological models is  
539 negligible (e.g., Viney et al., 2009b; Chiew, 2010; Petheram et al., 2012). However, in the current  
540 study, XAJ produced distinct results from the other models (see Figure 5 results with 1 donor  
541 catchment), which suggests that the performance of regionalization methods is affected by the choice  
542 of hydrological models even with one donor catchment.

543

## 544 5.2.2 Assessment over hydrological models

545

546 **Although** we claimed that the methods with **the** output averaging option (SP-out and Phy-out)  
547 produced better performance than **the** other methods, it is **difficult** to **determine** the most appropriate  
548 method between **the** spatial proximity (SP-out) and physical similarity (Phy-out) methods (also valid  
549 for excluding the influence **on the** hydrological model performance of calibration and verification, see  
550 Appendix B). This is consistent with the evaluation by using one hydrological model (monthly  
551 WASMOD) in the same area by Yang et al. (2018). According to the explanation from Oudin et al.  
552 (2008), it is not possible to decide which approach (SP-out or Phy-out) is the most appropriate one  
553 when the streaming network density is lower than 60 stations per 100,000 km<sup>2</sup>. As we used four  
554 hydrological models at different complexity levels, this result additionally confirmed that this  
555 assertion is independent of the selection of hydrological models.

556

557 Investigating the model preference for regionalization methods from different aspects, XAJ should be  
558 preferred when the evaluation is more focused on high flow, while WASMOD should be considered  
559 for **low-flow** analysis. This result is consistent with the model performance for gauged catchments (see  
560 Figure 4 and Table 6). This result tends to **support** the claim that there is no incentive to prefer a  
561 parsimonious hydrological model for regionalization **studies** rather than a model with adequate  
562 complexity (Arsenault et al., 2015; Poissant et al., 2017). However, hydrological models with fewer  
563 parameters are recommended when no preknowledge about the regionalization performance is  
564 available since the performance difference between **the** regionalization methods is relatively smaller.  
565 For the regression method, the model with more parameters works worse, probably due to the stronger  
566 interaction influence when increasing the number of parameters (e.g., Perrin et al., 2003; Poissant et al.,  
567 2017). Another **limitation of** the regression method **is** that not all the functions for **the** model  
568 parameters follow the linear assumption (e.g., Blöschl, 2005) and poor performance results from the  
569 accumulated errors.

570

### 571 5.2.3 Assessment in different climate regions

572 According to both the NSE and NSElog results, SP-out and Phy-out perform best for all the climate  
573 regions. Therefore, it seems reasonable to conclude that the selection of the climatic region has no  
574 large effect on the ranking of regionalization methods. However, the average regionalization  
575 performance in the oceanic climate region is substantially better and varies within a smaller range than  
576 in the other two cold regions. This indicates that the uncertainty in the selection of regionalization  
577 methods is larger in cold and dry regions than in warm and wet regions (see Figure 2). Due to the  
578 limited number of catchments in the oceanic climate and polar tundra climate regions, further  
579 comprehensive studies are needed to conclude the preferences of hydrological models and  
580 regionalization methods over various regions.

581

## 582 **6. Conclusions**

583 The main aim of this study was to investigate how different combinations of regionalization methods,  
584 hydrological models and climate conditions will influence the overall performance of hydrological  
585 simulations in ungauged basins. We assessed the performance of four hydrological models and five  
586 regionalization schemes (a) under stationary climate conditions to test how the performance of the  
587 regionalization methods depends on the choice of hydrological models, (b) under different climate  
588 conditions to assess the stability in performance of the hydrological models and regionalization  
589 methods as climate changes, and (c) in different climate regions to test how the performances of the  
590 simulations vary between these regions. The study was performed using data from 86 catchments in  
591 Norway, covering three climatic groups according to the Köppen-Geiger classification.

592

593 In this study, we found that for all the hydrological models, the distance-based approaches with the  
594 output averaging option (SP-out and Phy-out) always outperformed the other tested methods,  
595 especially for the low-flow estimation. Second, the difference in performance between the output and



596 parameter averaging options is not stable and positively increases with the number of parameters for  
597 the hydrological models. From our study, the performance difference between these options is the  
598 largest for XAJ and the smallest for GR4J. Third, the performance difference among the  
599 regionalization methods was smaller for models with fewer parameters (GR4J and WASMOD)  
600 compared to that of the models with more tunable parameters (HBV and XAJ). Regarding the model  
601 influence on regionalization performance, XAJ is recommended as the best-performing model  
602 according to the evaluations by NSE and F values, whereas NSElog recommends WASMOD as the  
603 best through the evaluation. Furthermore, clear differences in general were displayed for three climatic  
604 regions, and oceanic climatic regions provided the best performance and smallest variance over the  
605 regionalization methods and hydrological models. Moreover, the difference in hydrological model  
606 performance seems smaller among the regionalization methods than among the climate regions. From  
607 calibration to verification periods, the general performance for the regionalization methods did not  
608 show large degradations.

609

610 Although this study produced some solid conclusions that were not available before, there are some  
611 limitations of the current study. Compared with the general evaluation of hydrograph fit and water  
612 balance, assessment with emphasis on low flow showed more contrasting results, which requires  
613 closer attention in future work. In addition, studies with more different hydrological models are  
614 needed to show the influence of hydrological model selection on regionalization performance.  
615 Moreover, studies with more contrast in climate conditions are recommended to investigate the  
616 transferability of conclusions across climate regions and climate changing conditions, which is  
617 essential for future prediction.

618

619 **Acknowledgments:** This work is supported by the Research Council of Norway (FRINATEK Project  
620 274310), Research and Development Funding (Project number 80203) of the Norwegian Water  
621 Resources and Energy Directorate (NVE), and the China Scholarship Council. We would like to thank  
622 the NVE for providing the data for this study. We are thankful to the reviewers whose insightful and  
623 constructive comments have led to a significant improvement in the quality of the paper.

## **Reference**

626 Allen, R. G., Pereira, L. S., Raes, D., and Smith, M., 1998. Crop evapotranspiration, guidelines  
627 for computing crop water requirements, Irrig. and Drain. Pap. 56. U.N. Food and Agric.  
628 Organ., Rome.

629 Arsenault, R., Brissette, F.P., 2014. Continuous streamflow prediction in ungauged basins: The  
630 effects of equifinality and parameter set selection on uncertainty in regionalization  
631 approaches. *Water Resour. Res.* 50, 6135–6153. <https://doi.org/10.1002/2013WR014898>

632 Arsenault, R., Poissant, D., Brissette, F., 2015. Parameter dimensionality reduction of a  
633 conceptual model for streamflow prediction in Canadian, snowmelt dominated ungauged  
634 basins. *Adv. Water Resour.* 85, 27–44. <https://doi.org/10.1016/j.advwatres.2015.08.014>

635 Bao, Z., Zhang, J., Liu, J., Fu, G., Wang, G., He, R., Yan, X., Jin, J., Liu, H., 2012. Comparison  
636 of regionalization approaches based on regression and similarity for predictions in  
637 ungauged catchments under multiple hydro-climatic conditions. *J. Hydrol.* 466–467, 37–46.  
638 <https://doi.org/10.1016/j.jhydrol.2012.07.048>

639 Bárdossy, A., 2007. Calibration of hydrological model parameters for ungauged catchments.  
640 *Hydrol. Earth Syst. Sci. Discuss.* 3, 1105–1124. <https://doi.org/10.5194/hessd-3-1105-2006>

641 Beck, H.E., Zimmermann, N.E., McVicar, T.R., Vergopolan, N., Berg, A., Wood, E.F., 2018.  
642 Present and future köppen-geiger climate classification maps at 1-km resolution. *Sci. Data*  
643 5, 1–12. <https://doi.org/10.1038/sdata.2018.214>

644 Blöschl, G., 2005. Rainfall–runoff modelling of ungauged catchments. In: Anderson, M.G. (Ed.),  
645 *Encyclopedia of Hydrological Sciences*. John Wiley & Sons, Chichester, pp. 2061–2080.

646 Blöschl, G., Montanari, A., 2010. Climate change impacts-throwing the dice? *Hydrol. Process* 24,  
647 374–381. <https://doi.org/10.1002/hyp.7574>

648 Burn, D.H., Boorman, D.B., 1993. Estimation of hydrological parameters at ungauged  
649 catchments. *J. Hydrol.* 143, 429–454. [https://doi.org/10.1016/0022-1694\(93\)90203-L](https://doi.org/10.1016/0022-1694(93)90203-L)

650 Broderick C., 2016. Transferability of hydrological models and ensemble averaging methods  
651 between contrasting climatic periods Ciaran. *Water Resour. Res.* 52, 8343–8373.  
652 <https://doi.org/10.1002/2016WR018850>.

653 Chiew, F.H.S., 2010. Lumped conceptual rainfall-runoff models and simple water balance  
654 methods: Overview and applications in ungauged and data limited regions. *Geogr. Compass*  
655 4, 206–225. <https://doi.org/10.1111/j.1749-8198.2009.00318.x>

656 Coron, L., Andre, V., Perrin, C., Lerat, J., Vaze, J., Bourqui, M., Hendrickx, F., 2012. Crash  
657 testing hydrological models in contrasted climate conditions : An experiment on 216  
658 Australian catchments 48, 1–17. <https://doi.org/10.1029/2011WR011721>

659 Coron, L., Andréassian, V., Perrin, C., Bourqui, M., Hendrickx, F., 2014. On the lack of  
660 robustness of hydrologic models regarding water balance simulation: A diagnostic  
661 approach applied to three models of increasing complexity on 20 mountainous catchments.  
662 *Hydrol. Earth Syst. Sci.* 18, 727–746. <https://doi.org/10.5194/hess-18-727-2014>

663 Egbuniwe, N., Todd, D.K., 1976. Application of the Stanford Watershed Model To Nigerian  
664 Watersheds. *JAWRA J. Am. Water Resour. Assoc.* 12, 449–460.  
665 <https://doi.org/10.1111/j.1752-1688.1976.tb02710.x>

666 He, Y., Bárdossy, A., Zehe, E., 2011. A review of regionalization for continuous streamflow  
667 simulation. *Hydrol. Earth Syst. Sci.* 15, 3539–3553. [https://doi.org/10.5194/hess-15-3539-](https://doi.org/10.5194/hess-15-3539-2011)  
668 [2011](https://doi.org/10.5194/hess-15-3539-2011).

669 Hargreaves, G. H.: 1975, ‘Moisture Availability and Crop Production’, *TRANSACTION of the*  
670 *ASAE* 18, 980–984.

671 Hargreaves, G. H., Samani, Z. A., (1985). Reference crop evapotranspiration from temperature.  
672 *Appl. Eng. Agric.*, 1(2), 96–99.

673 Holländer, H.M., Blume, T., Bormann, H., Buytaert, W., Chirico, G.B., Exbrayat, J.F.,  
674 Gustafsson, D., Hölzel, H., Kraft, P., Stamm, C., Stoll, S., Blöschl, G., Flühler, H., 2009.  
675 Comparative predictions of discharge from an artificial catchment (Chicken Creek) using  
676 sparse data. *Hydrol. Earth Syst. Sci.* 13, 2069–2094. [https://doi.org/10.5194/hess-13-2069-](https://doi.org/10.5194/hess-13-2069-2009)  
677 [2009](https://doi.org/10.5194/hess-13-2069-2009)

678 Hrachowitz, M., Savenije, H.H.G., Blöschl, G., McDonnell, J.J., Sivapalan, M., Pomeroy, J.W.,  
679 Arheimer, B., Blume, T., Clark, M.P., Ehret, U., Fenicia, F., Freer, J.E., Gelfan, A., Gupta,  
680 H.V., Hughes, D. a., Hut, R.W., Montanari, A., Pande, S., Tetzlaff, D., Troch, P.A.,  
681 Uhlenbrook, S., Wagener, T., Winsemius, H.C., Woods, R. a., Zehe, E., Cudennec, C.,  
682 2013. A decade of Predictions in Ungauged Basins (PUB)—a review. *Hydrol. Sci. J.* 58,  
683 1198–1255. <https://doi.org/10.1080/02626667.2013.803183>

684 Hublart, P., Ruelland, D., García De Cortázar Atauri, I., Ibacache, A., 2015. Reliability of a  
685 conceptual hydrological model in a semi-arid Andean catchment facing water-use changes,  
686 in: *IAHS-AISH Proceedings and Reports*. pp. 203–209. [https://doi.org/10.5194/piahs-371-](https://doi.org/10.5194/piahs-371-203-2015)  
687 [203-2015](https://doi.org/10.5194/piahs-371-203-2015)

688 IPCC, 2014. *Climate Change 2014: Synthesis Report. Contribution of Working Groups I, II and*  
689 *III to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*. In:  
690 *Core Writing Team, R.K. Pachauri and L.A. Meyer (Eds.). IPCC, Geneva, Switzerland, 151*  
691 *pp.*

692 Jansson, A, Tveito, O E, Pirinen, P, & Scharling, M., 2007. *NORDGRID - a preliminary*  
693 *investigation on the potential for creation of a joint Nordic gridded climate dataset*. met.no  
694 *Report 03/2007.*

695 Jin, X., Xu, C. yu, Zhang, Q., Chen, Y.D., 2009. Regionalization study of a conceptual  
696 hydrological model in Dongjiang basin, south China. *Quat. Int.* 208, 129–137.  
697 <https://doi.org/10.1016/j.quaint.2008.08.006>

698 Klemeš, V., 1986. Operational testing of hydrological simulation models. *Hydrol. Sci. J.* ISSN  
699 6667. [https://doi.org/10.1080/02626668609491024.](https://doi.org/10.1080/02626668609491024)

700 Kottek, M., Grieser, J., Beck, C., Rudolf, B., Rubel, F., 2006. World Map of the Köppen-Geiger  
701 climate classification updated. *Meteorologische Zeitschrift*, 15, 259–263.

702 Lagarias, J.C., Reeds, J.A., Wright, M.H., Wright, P.E., 1998. Convergence Properties of the  
703 Nelder--Mead Simplex Method in Low Dimensions. *SIAM J. Optim.* 9, 112–147.  
704 <https://doi.org/10.1137/S1052623496303470>

705 Li, C.Z., Zhang, L., Wang, H., Zhang, Y.Q., Yu, F.L., Yan, D.H., 2012. The transferability of  
706 hydrological models under nonstationary climatic conditions 1239–1254.  
707 <https://doi.org/10.5194/hess-16-1239-2012>

708 Li, F., Zhang, Y., Xu, Z., Liu, C., Zhou, Y., Liu, W., 2014. Runoff predictions in ungauged  
709 catchments in southeast Tibetan Plateau. *J. Hydrol.* 511, 28–38.  
710 <https://doi.org/10.1016/j.jhydrol.2014.01.014>

711 Li, H., Zhang, Y., 2017. Regionalising rainfall-runoff modelling for predicting daily runoff:  
712 Comparing gridded spatial proximity and gridded integrated similarity approaches against  
713 their lumped counterparts. *J. Hydrol.* 550. <https://doi.org/10.1016/j.jhydrol.2017.05.015>

714 Li, H., Zhang, Y., Chiew, F.H.S., Xu, S., 2009. Predicting runoff in ungauged catchments by  
715 using Xinanjiang model with MODIS leaf area index. *J. Hydrol.* 370, 155–162.  
716 <https://doi.org/10.1016/j.jhydrol.2009.03.003>

717 Li, L., Diallo, I., Xu, C.Y., Stordal, F., 2015. Hydrological projections under climate change in  
718 the near future by RegCM4 in Southern Africa using a large-scale hydrological model. *J.*  
719 *Hydrol.* 528, 1–16. <https://doi.org/10.1016/j.jhydrol.2015.05.028>

720 Li, L., Ngongondo, C.S., Xu, C.-Y., Gong, L., 2013. Comparison of the global TRMM and WFD  
721 precipitation datasets in driving a large-scale hydrological model in southern Africa. *Hydrol.*  
722 *Res.* 44, 770. <https://doi.org/10.2166/nh.2012.175>

723 Lin K., Liu P., He Y., Guo S., 2014. Multi-site evaluation to reduce parameter uncertainty in a  
724 conceptual hydrological modeling within the GLUE framework. *J. HYDROINFORM.* 16  
725 (1), 60–73. doi: <https://doi.org/10.2166/hydro.2013.204>

726 Magette, W.L., Shanholtz, V.O., Carr, J.C., 1976. Estimating selected parameters for the  
727 Kentucky Watershed Model from watershed characteristics. *Water Resour. Res.* 12, 472–  
728 476. <https://doi.org/10.1029/WR012i003p00472>

729 McIntyre, N., Lee, H., Wheeler, H., Young, A., Wagener, T., 2005. Ensemble predictions of  
730 runoff in ungauged catchments. *Water Resour. Res.* 41, 1–14.  
731 <https://doi.org/10.1029/2005WR004289>

732 Merz, R., Blöschl, G., 2004. Regionalization of catchment model parameters. *J. Hydrol.* 287, 95–  
733 123. <https://doi.org/10.1016/j.jhydrol.2003.09.028>

734 Mohr, M., 2009. Comparison of Version 1.1 and 1.0 of gridded temperature and precipitation  
735 data for Norway. *met.no Note 19/2009*.

736 Moriasi, D.N., Arnold, J.G., VanLiew, M.W., Bingner, R.L., Harmel, R.D., Veith, T.L., 2007.  
737 Evaluation guidelines for systematic quantification of accuracy in watershed simulations.  
738 *Trans. ASABE* 50(3), 885–900.

739 Muller-Wohlfeil, D.-I., XU, C.-Y., Lversen, H.L., 2003. Estimation of Monthly River Discharge  
740 from Danish Catchments Background and Objective of the Study. *Nord. Hydrol.* 34, 295–  
741 320.

742 Nash, E., Sutcliffe, V., 1970. RIVER FLOW FORECASTING THROUGH CONCEPTUAL  
743 MODELS PART I- A DISCUSSION OF PRINCIPLES. *J. Hydrol.* 10, 282–290.

744 Oudin, L., Andréassian, V., Perrin, C., Michel, C., Le Moine, N., 2008. Spatial proximity,  
745 physical similarity, regression and ungauged catchments: A comparison of regionalization  
746 approaches based on 913 French catchments. *Water Resour. Res.* 44, 1–15.  
747 <https://doi.org/10.1029/2007WR006240>

748 Parajka, J., Blöschl, G., Merz, R., 2007. Regional calibration of catchment models: Potential for  
749 ungauged catchments. *Water Resour. Res.* 43. <https://doi.org/10.1029/2006WR005271>

750 Parajka, J., Merz, R., Blöschl, G., 2005. A comparison of regionalization methods for catchment  
751 model parameters. *Hydrol. Earth Syst. Sci. Discuss.* 2, 509–542.  
752 <https://doi.org/10.5194/hessd-2-509-2005>

753 Parajka, J., Viglione, A., Rogger, M., Salinas, J.L., Sivapalan, M., Blöschl, G., 2013.  
754 Comparative assessment of predictions in ungauged basins-Part 1: Runoff-hydrograph  
755 studies. *Hydrol. Earth Syst. Sci.* 17, 1783–1795. <https://doi.org/10.5194/hess-17-1783-2013>

756 Peel, M.C., Finlayson, B.L., McMahon, T.A., 2007. Updated world map of the Köppen-Geiger  
757 climate classification. *Hydrol. Earth Syst. Sci.* 11, 1633–1644.

758 Perrin, C., Michel, C., Andreassian, V., 2001. Does a large number of parameters enhance model  
759 performance? Comparative assessment of common catchment model structures on 429  
760 catchments. *J. Hydrol.* 242, 275–301.

761 Perrin, C., Michel, C., Andréassian, V., 2003. Improvement of a parsimonious model for  
762 streamflow simulation. *J. Hydrol.* 279, 275–289. [https://doi.org/10.1016/S0022-](https://doi.org/10.1016/S0022-1694(03)00225-7)  
763 [1694\(03\)00225-7](https://doi.org/10.1016/S0022-1694(03)00225-7)

764 Petheram, C., Rustomji, P., Chiew, F.H.S., Vleeshouwer, J., 2012. Rainfall-runoff modelling in  
765 northern Australia: A guide to modelling strategies in the tropics. *J. Hydrol.* 462–463, 28–  
766 41. <https://doi.org/10.1016/j.jhydrol.2011.12.046>

767 Poissant, D., Arsenault, R., Brissette, F., 2017. Impact of parameter set dimensionality and  
768 calibration procedures on streamflow prediction at ungauged catchments. *J. Hydrol. Reg.*  
769 *Stud.* 12, 220–237. <https://doi.org/10.1016/j.ejrh.2017.05.005>

770 Pool, S., Viviroli, D., Seibert, J., 2017. Prediction of hydrographs and flow-duration curves in  
771 almost ungauged catchments: Which runoff measurements are most informative for model  
772 calibration? *J. Hydrol.* 554, 613–622. <https://doi.org/10.1016/j.jhydrol.2017.09.037>

773 Razavi, T., Coulibaly, P., 2013. Streamflow Prediction in Ungauged Basins: Review of  
774 Regionalization Methods. *J. Hydrol. Eng.* 18, 958–975.  
775 [https://doi.org/10.1061/\(ASCE\)HE.1943-5584.0000690](https://doi.org/10.1061/(ASCE)HE.1943-5584.0000690)

776 Reichl, J.P.C., Western, A.W., McIntyre, N.R., Chiew, F.H.S., 2009. Optimization of a similarity  
777 measure for estimating ungauged streamflow. *Water Resour. Res.* 45, 1–15.  
778 <https://doi.org/10.1029/2008WR007248>

779 Rojas-Serna, C., Lebecherel, L., Perrin, C., Andréassian, V., Oudin, L., 2016. How should a  
780 rainfall-runoff model be parameterized in an almost ungauged catchment? A methodology  
781 tested on 609 catchments. *Water Resour. Res.* 1–24.  
782 <https://doi.org/10.1002/2016WR018704>.Received

783 Salinas, J.L., Laaha, G., Rogger, M., Parajka, J., Viglione, A., Sivapalan, M., Blöschl, G., 2013.  
784 Comparative assessment of predictions in ungauged basins-Part 2: Flood and low flow  
785 studies. *Hydrol. Earth Syst. Sci.* 17, 2637–2652. <https://doi.org/10.5194/hess-17-2637-2013>

786 Samuel, J., Coulibaly, P., Metcalfe, R.A., 2011. Estimation of Continuous Streamflow in Ontario  
787 Ungauged Basins: Comparison of Regionalization Methods. *J. Hydrol. Eng.* 16, 447–459.  
788 [https://doi.org/10.1061/\(ASCE\)HE.1943-5584.0000338](https://doi.org/10.1061/(ASCE)HE.1943-5584.0000338)

789 Seibert, J., Beven, K.J., 2009. Gauging the ungauged basin : how many discharge measurements  
790 are needed? *Hydrol. Earth Syst. Sci.* 13, 883–892. [https://doi.org/10.5194/hessd-6-2275-](https://doi.org/10.5194/hessd-6-2275-2009)  
791 2009

792 Seibert, J., Vis, M.J.P., 2012. Teaching hydrological modeling with a user-friendly catchment-  
793 runoff-model software package. *Hydrol. Earth Syst. Sci.* 16, 3315–3325.  
794 <https://doi.org/10.5194/hess-16-3315-2012>

795 Shuttleworth, W.J., 1993. Evaporation. In: Maidment, D.R. (Ed.), *Handbook of Hydrology*.  
796 McGraw-Hill, New York, pp. 4.1–4.53.

797 Sivapalan, M., Takeuchi, K., Franks, S.W., Gupta, V.K., Karambiri, H., Lakshmi, V., Liang, X.,  
798 McDonnell, J.J., Mendiondo, E.M., O’Connell, P.E., Oki, T., Pomeroy, J.W., Schertzer, D.,  
799 Uhlenbrook, S., Zehe, E., 2003. IAHS Decade on Predictions in Ungauged Basins (PUB),  
800 2003–2012: Shaping an exciting future for the hydrological sciences. *Hydrol. Sci. J.* 48,  
801 857–880. <https://doi.org/10.1623/hysj.48.6.857.51421>

802 Tveito, O E, Bjørndal, I, Skjelvåg, A O, & Aune, B., 2005. A GIS-based agro-ecological decision  
803 system based on gridded climatology. *Meteorol. Appl.* 12(1).

804 Valéry A. 2010. Modélisation précipitations–débit sous influence nivale.Élaboration d’un  
805 module neige et évaluation sur 380 bassins versants. Agro Paris Tech: Paris, France.

806 Vandewiele, G.L., Xu, C.Y., Huybrechts, W., 1991. Regionalization of physically-based water  
807 balance models in Belgium. Application to ungauged catchments. *Water Resour. Manag.* 5,  
808 199–208. <https://doi.org/10.1007/BF00421989>

809 Vaze, J., Post, D.A., Chiew, F.H.S., Perraud, J., Viney, N.R., Teng, J., 2010. Climate non-  
810 stationarity – Validity of calibrated rainfall – runoff models for use in climate change  
811 studies. *J. Hydrol.* 394, 447–457. <https://doi.org/10.1016/j.jhydrol.2010.09.018>



812 Viglione, A., Parajka, J., Rogger, M., Salinas, J.L., Laaha, G., Sivapalan, M., Blöschl, G., 2013.  
813 Comparative assessment of predictions in ungauged basins - Part 3: Runoff signatures in  
814 Austria. *Hydrol. Earth Syst. Sci.* 17, 2263–2279. <https://doi.org/10.5194/hess-17-2263-2013>

815 Viney, N.R., Perraud, J., Vaze, J., Chiew, F.H.S., Post, D.A., Yang, A., 2009a. The usefulness of  
816 bias constraints in model calibration for regionalization to ungauged catchments. 18<sup>th</sup>  
817 World IMACS/MODSIM Congr. Cairns, Aust. 3421– 3427.

818 Viney, N.R., Vaze, J., Chiew, F.H.S., Perraud, J.M., Post, D.A., Teng, J., 2009b. Comparison of  
819 multi-model and multi-donor ensembles for regionalization of runoff generation using five  
820 lumped rainfall–runoff models. In: MODSIM 2009 International Congress on Modelling  
821 and Simulation. MSSANZ, Cairns, Australia, pp. 3428–3434.

822 Vormoor, K., Lawrence, D., Schlichting, L., Wilson, D., Kwok, W., 2016. Evidence for changes  
823 in the magnitude and frequency of observed rainfall vs . snowmelt driven floods in Norway.  
824 *J. Hydrol.* 538, 33–48. <https://doi.org/10.1016/j.jhydrol.2016.03.066>

825 Widén-Nilsson, E., Halldin, S., Xu, C. y., 2007. Global water-balance modelling with  
826 WASMOD-M: Parameter estimation and regionalization . *J. Hydrol.* 340, 105–118.  
827 <https://doi.org/10.1016/j.jhydrol.2007.04.002>

828 Xu, C., 1999. Operational testing of a water balance model for predicting climate change impacts.  
829 *Agric. Forest Meteorol.* 23, 95–304.

830 Xu, C., Halldin, S., 1997. The Effect of Climate Change on River Flow and Snow Cover in the  
831 NOPEX Area Simulated by a Simple Water Balance Model. *Nordic Hydrol.* 28 (4/5), 273–  
832 282.

833 Xu, C., Singh, V.P., 2002. Cross Comparison of Empirical Equations for Calculating Potential  
834 Evapotranspiration with Data from Switzerland. *Water Resour. Manag.* 16, 197–219.

835 Xu, C.Y., 2003. Testing the transferability of regression equations derived from small sub-  
836 catchments to a large area in central Sweden. *Hydrol. Earth Syst. Sci.* 7, 317–324.  
837 <https://doi.org/10.5194/hess-7-317-2003>

838 Xu, C.Y., 1999. Estimation of parameters of a conceptual water balance model for ungauged  
839 catchments. *Water Resour. Manag.* 13, 353–368. <https://doi.org/10.1023/A:1008191517801>

840 Yang, X., Magnusson, J., Rizzi, J., Xu, C., 2018. Runoff prediction in ungauged catchments in  
841 Norway: comparison of regionalization approaches. *Hydrol. Res.* 49(2), 487-505.  
842 <https://doi.org/10.2166/nh.2017.071>

843 Yang, X., Magnusson, J., Xu, C.-Y., 2019. Transferability of regionalization methods under  
844 changing climate. *J. Hydrol.* 568, 67-81. <https://doi.org/10.1016/j.jhydrol.2018.10.030>

845 Young, A.R., 2006. Stream flow simulation within UK ungauged catchments using a daily  
846 rainfall-runoff model. *J. Hydrol.* 320, 155–172.  
847 <https://doi.org/10.1016/j.jhydrol.2005.07.017>

848 Zhang, Y., Chiew, F., 2009. Evaluation of Regionalization Methods for Predicting Runoff in  
849 Ungauged Catchments in Southeast Australia. 18<sup>th</sup> World IMACS/MODSIM Congr. Cairns,  
850 Aust. 3442–3448.

851 Zhang, Y., Vaze, J., Chiew, F.H.S., Teng, J., Li, M., 2014. Predicting hydrological signatures in  
852 ungauged catchments using spatial interpolation , index model , and rainfall – runoff  
853 modelling. *J. Hydrol.* 517, 936–948. <https://doi.org/10.1016/j.jhydrol.2014.06.032>

854 Zhang, Y., Zheng, H., Chiew, F.H.S., Arancibia, J.P., Zhou, X., 2016. Evaluating Regional and  
855 Global Hydrological Models against Streamflow and Evapotranspiration Measurements. *J.*  
856 *Hydrometeor.* 17, 995–1010. <https://doi.org/10.1175/JHM-D-15-0107.1>

857 Zhang, Z.X., Chen, X., Xu, C-Y, Yuan, L.F., Yong, B., Yan, S.F., 2011a. Evaluating the non-  
858 stationary relationship between Precipitation and Streamflow in Nine Major Basins of  
859 China during the past 50 years. *J. Hydrol.* 409, 81-93.

860 Zhao R-J. (1992) Xinanjiang model applied in China. *J. Hydrol.* 135(2), 371–381.

861 Zhao R-J, Zuang Y, Fang L, Liu X, Zhang Q (1980). The Xinanjiang Model. In: *Hydrological*  
862 *Forecasting*. IAHS Press, Wallingford, pp. 351–356.

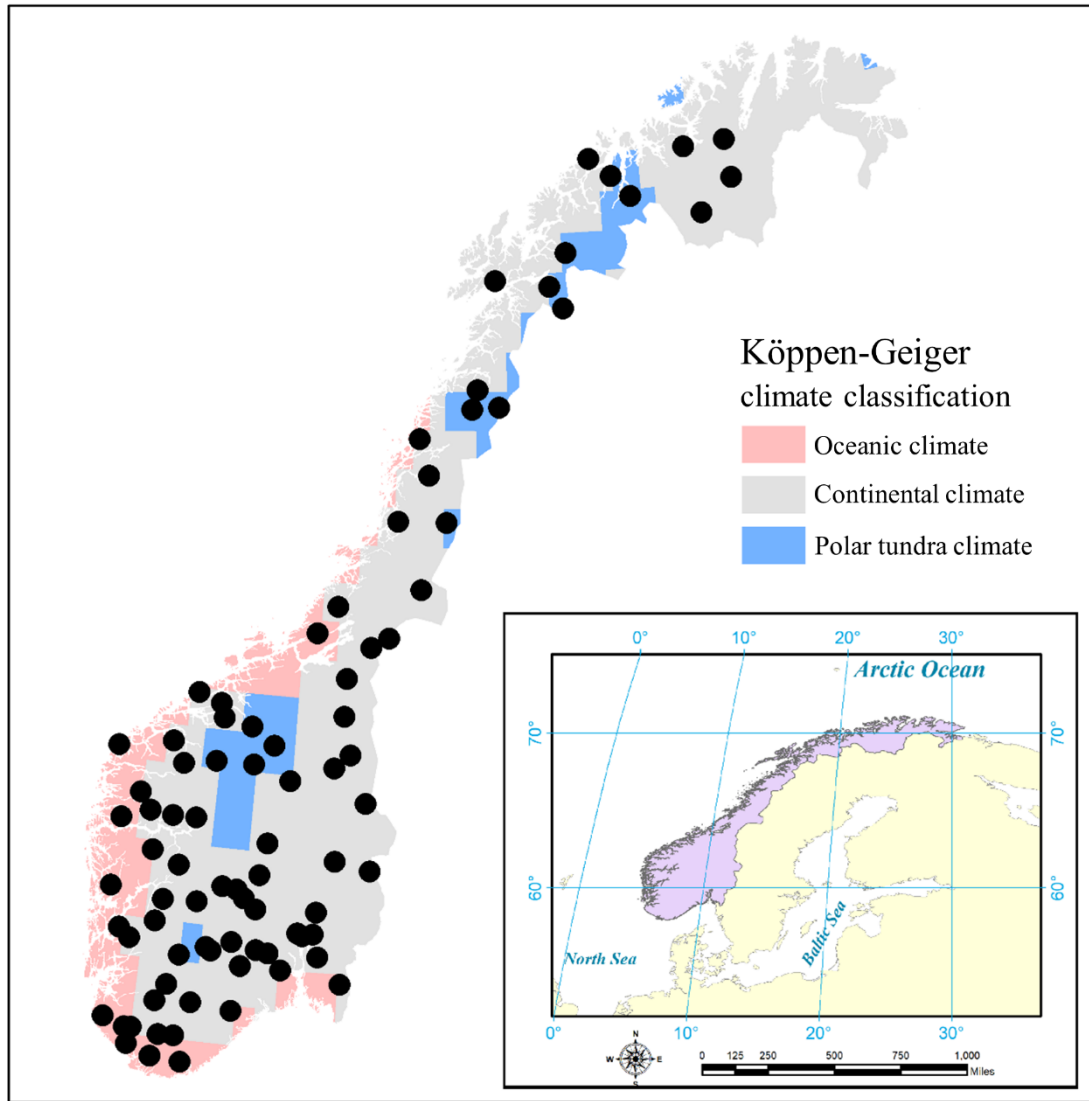


Fig. 1

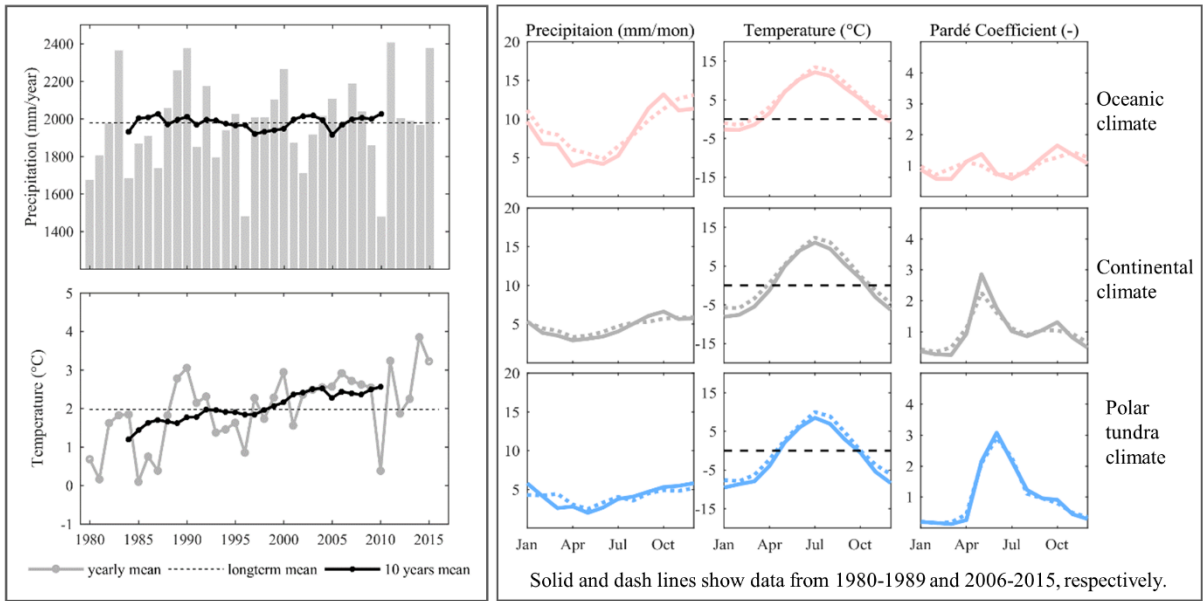
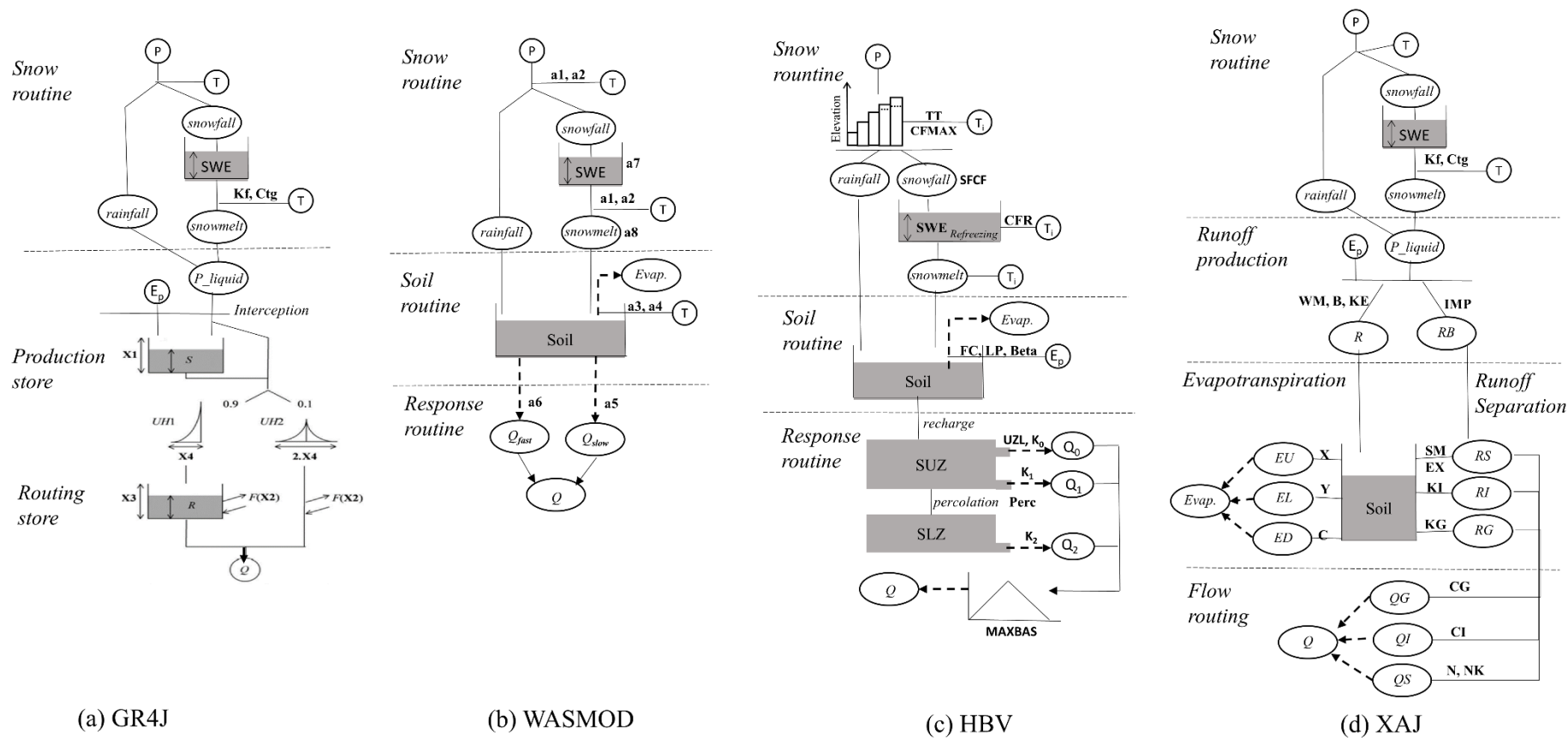


Fig. 2



P: Precipitation (mm/day) T: Temperature (°C) SWE: Snow Water Equivalent (mm) E: Evapotranspiration (mm/day) Q: Discharge (m<sup>3</sup>/s)

Fig. 3

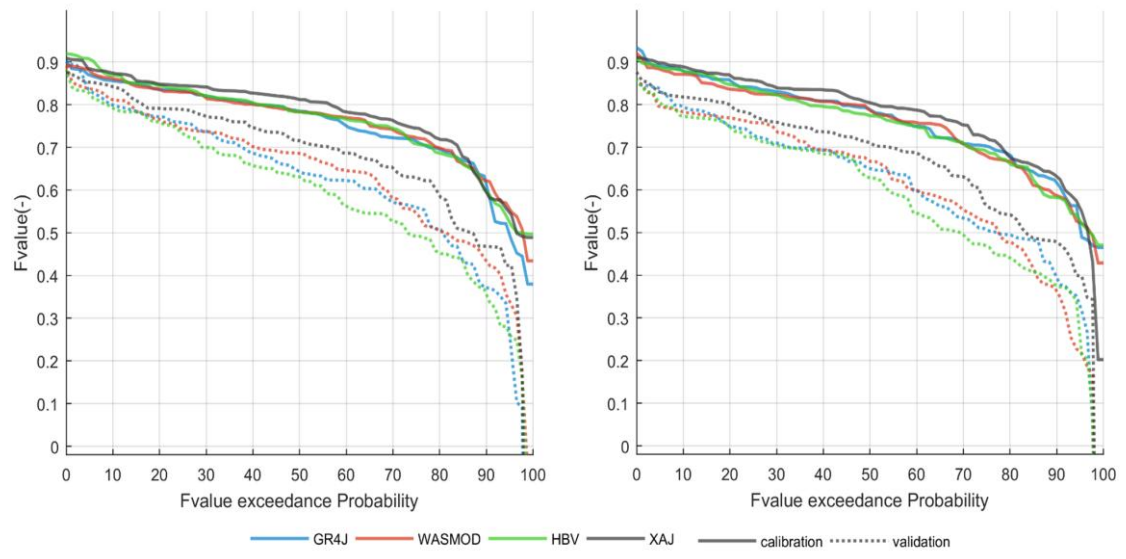


Fig. 4

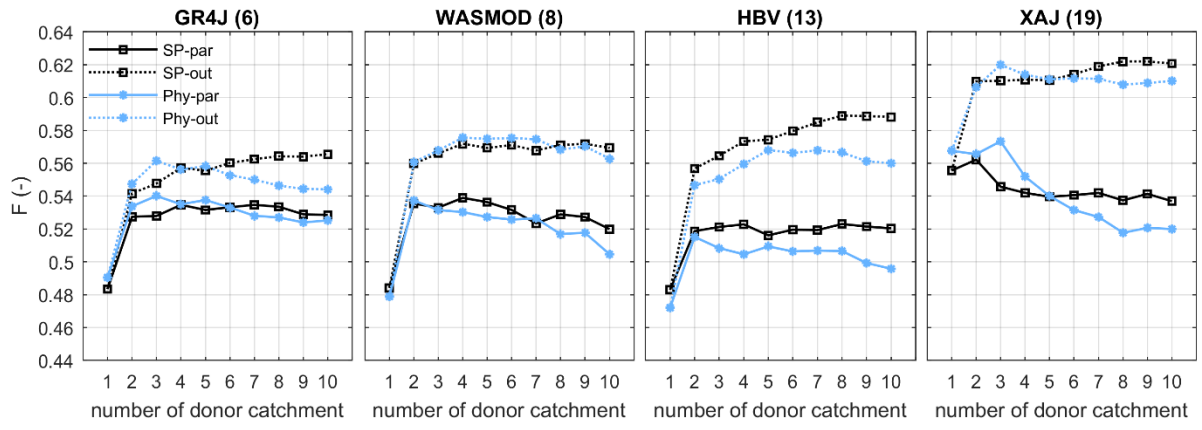


Fig. 5

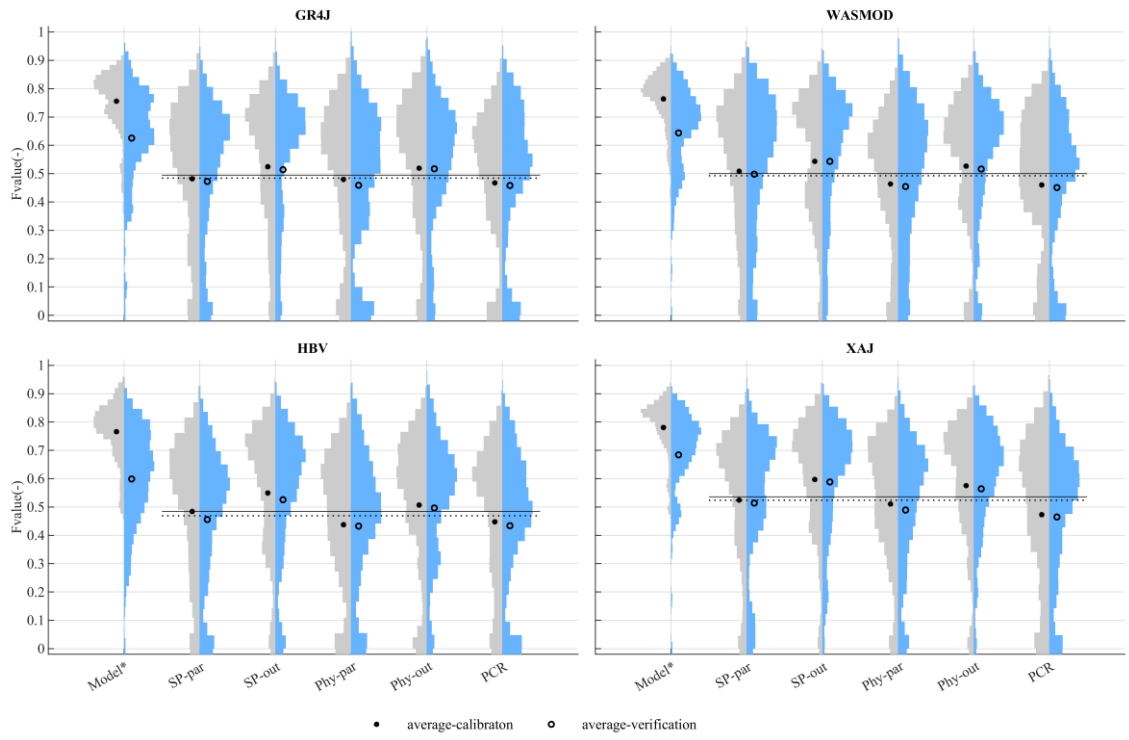


Fig. 6



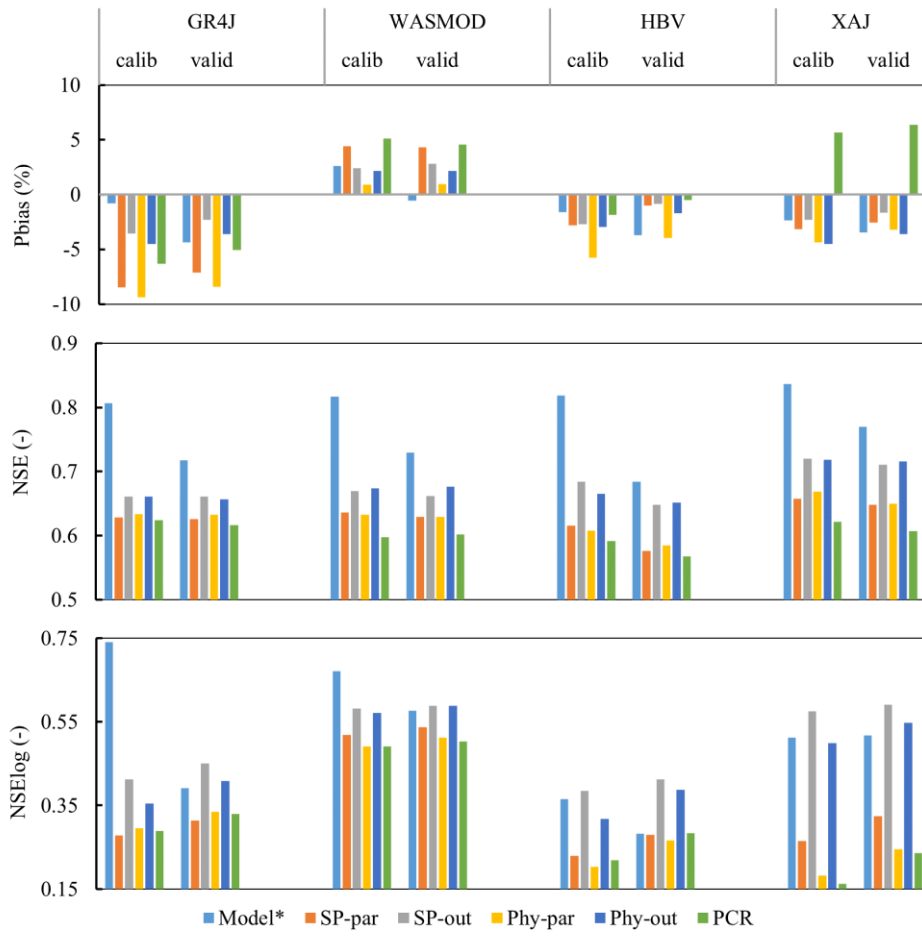


Fig. 7

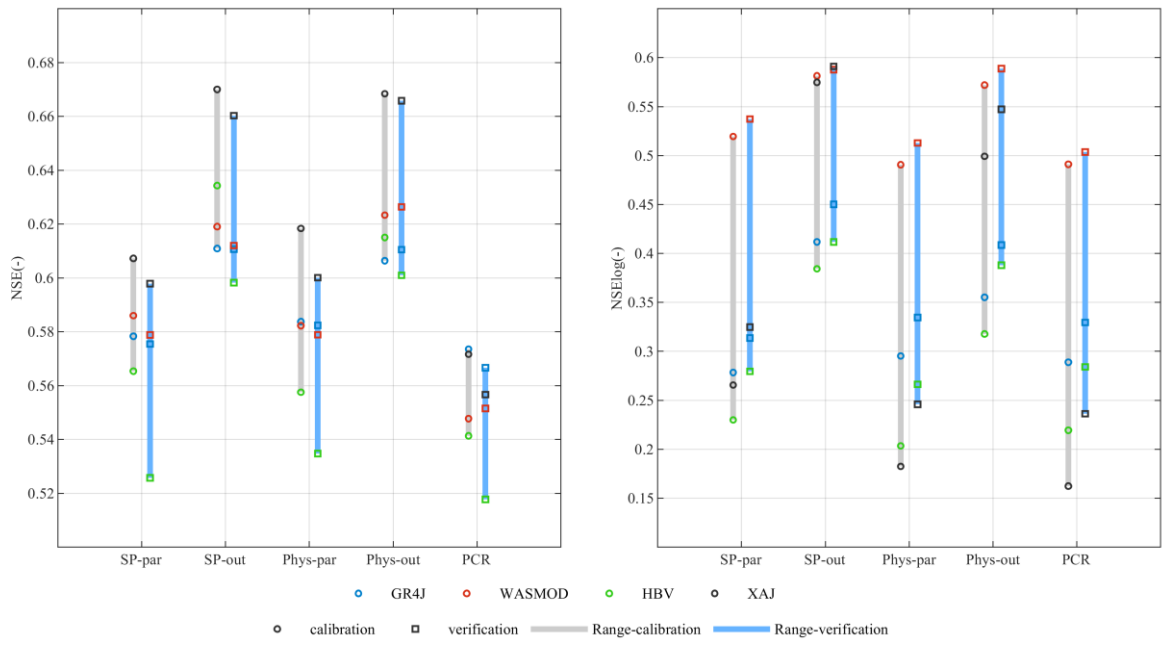


Fig. 8

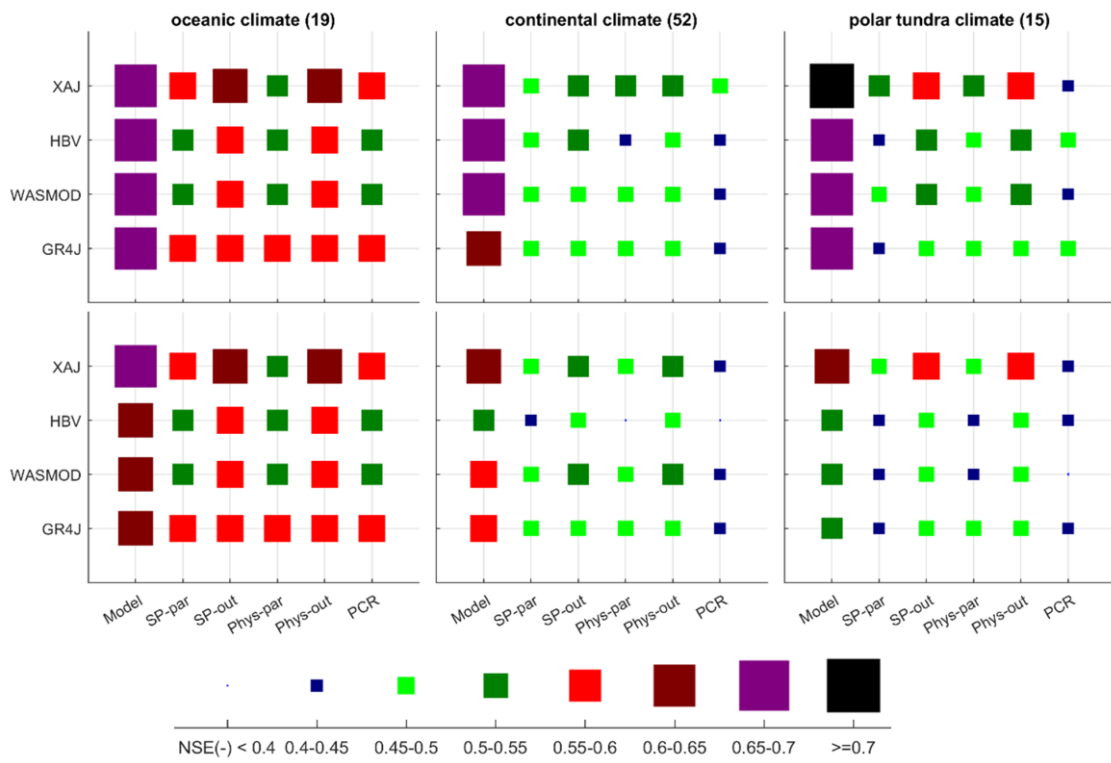


Fig. 9

**Figure captions:**

Fig.1 The location of the study catchments and the modified Köppen-Geiger climate classification.

Fig.2 The left panel shows yearly mean precipitation and temperature for the available data period, including a moving average with a sample window covering 10 years of data. The right panel shows the climatological distribution of precipitation, temperature and Pardé coefficient (i.e., ratio of the average monthly discharge to the mean annual runoff) using monthly data for the three climatological regions.

Fig.3 The structure of hydrological models tested in this study. The circles show the input variables, the ellipses present the process/output variable and the model parameters are marked with bold text. For detailed model equations, please refer to the references for the (a) GR4J model (Perrin et al., 2003; Valéry, 2010), (b) WASMOD model (Xu, 2003), (c) HBV model (Seibert and Vis, 2012), and (d) XAJ model (Lin et al., 2014).

Fig.4 The performance of hydrological models by split-sample test evaluated by the F value over 86 catchments. The left panel shows the results for calibration in 1980-1989 and verification in 2006-2015; the right panel displays the results of calibration in 2006-2015 and verification in 1980-1989.

Fig.5 Model performance versus number of donor catchments for the distance-based regionalization methods and four different models. The number of model parameters is given in the parenthesis next to the model name.

Fig.6 Split violin plots showing the distributions of F values for the five regionalization methods by each hydrological model during the calibration (gray color) and verification (blue color) periods. For each model and regionalization method, the solid black dots show the average performance for the calibration period, whereas the black circle shows the corresponding value for the verification period. The average performance of all regionalization methods for each hydrological model is shown as a solid line for the calibration period and as a dashed line for the verification period. The plot displays results from the 86 study catchments. The 'model' in the x-axis label shows the hydrological model performance in the calibration (gray color) and validation (blue color) periods.

Fig.7 Average performance for the different hydrological models and regionalization methods, given by Pbias, NSE and NSElog. Model\* is the result of model simulation performance in the calibration ('calib') and verification ('valid') periods.

Fig.8 Comparison of hydrological model performance over five regionalization methods in the calibration and verification periods. The bar shows the maximum difference between the hydrological models evaluated by the average NSE and NSElog values over 86 catchments.

Fig.9 The performance of the regionalization methods and hydrological models in different climatic regions. The size of the boxes is proportional to the average NSE value of the catchments within each climate group. The upper panel shows the results from the calibration period, and the lower panel shows the verification period. The number of catchments in each group is given in the title of each column. The 'Model' in the x-axis label shows the hydrological model performance for runoff simulation without regionalization.

Table 1. The average precipitation, temperature and runoff information for all climate groups

		Precipitation (mm/period)		Temperature (°C)		Runoff (mm/period)	
		calibration	validation	calibration	validation	calibration	validation
Oceanic climate	Year	2949	3211	4.1	5.2	2158	2342
	summer*	1411	1412	9.0	9.8	1197	1128
	winter	1508	1800	-0.7	0.5	961	1214
Continental climate	Year	1686	1750	0.8	2.3	1213	1250
	summer*	867	873	7.0	8.0	898	835
	winter	819	878	-5.3	-3.4	315	415
Polar tundra climate	Year	1633	1688	0.0	1.4	1187	1236
	summer*	817	819	6.1	7.1	942	908
	winter	816	869	-6.1	-4.3	245	328

\*Summer is from 1st of May to 31<sup>st</sup> of October.

Table 2. The statistical information about catchment descriptors used in regionalization methods

	Mean	Median	Minimum	Maximum
Area (km <sup>2</sup> )	340	145	3	5621
<u>Climate index</u>				
Mean annual precipitation (mm)	2255	1922	601	6008
Precipitation seasonality indices <sup>(1)</sup>	3.1	2.9	1.7	7.0
Mean annual temperature (°C)	2.7	2.5	-2.2	7.3
Temperature seasonality indices <sup>(2)</sup>	15.5	15.4	7.5	24.2
Aridity index <sup>(3)</sup>	0.1	0.1	0.0	0.4
<u>Terrain characteristics</u>				
Mean slope (°)	11	9	2	26
Mean elevation (m)	666	590	157	1472
<u>Land use</u>				
Artificial (%)	0.5	0.0	0.0	8.0
Agriculture (%)	4.1	1.1	0.0	57.6
Forest (%)	84.7	87.8	34.8	100.0
Wetland (%)	7.0	2.2	0.0	41.6
Waterbody (%)	3.7	2.9	0.0	15.1

(1) Precipitation seasonality indices: the ratio between the three consecutive wettest and driest months for each watershed.

(2) Temperature seasonality indices: the mean temperature of the hottest month minus the mean temperature of the coldest month in °C.

(3) Aridity index: the ratio between annual mean precipitation and potential evapotranspiration.

Table 3 Description of the calibrated model parameters in this study.

Parameter	Explanation	Reference	
<u><i>CemaNeige</i></u>			
$C_{TG}$	Ponderation coefficient	Valéry (2010)	
$K_f$	Degree-day factor		
<u><i>GR4J</i></u>			
X1	Production store maximal capacity	Perrin et al. (2003)	
X2	Catchment water exchange coefficient		
X3	One-day maximal capacity of the routing reservoir		
X4	<i>HUI</i> unit hydrograph time base		
<u><i>WASMOD</i></u>			
a1	Threshold temperature for rainfall and snowfall	Xu, (2003)	
a2	Threshold temperature for snowpack and snowmelt		
a3	Proportion parameter in potential evapotranspiration		
a4	Exponent parameter in actual evapotranspiration		
a5	Proportion coefficient of base flow		
a6	Proportion coefficient of fast flow		
a7	Coefficient for snowpack		
a8	Coefficient for snowmelt		
<u><i>HBV</i></u>			
TT	Threshold temperature	Seibert and Vis, (2012)	
CFMAX	Degree-day factor		
SFCF	Snowfall correction factor		
CFR	Refreezing coefficient		
FC	Field capacity		
LP	Threshold for reduction of evaporation		
Beta	Shape coefficient		
UZL	Threshold parameter for upper zone		
K0	Recession coefficient in upper zone		
K1	Recession coefficient in upper zone		
K2	Recession coefficient in lower zone		
Perc	Maximal flow from upper to lower box		
MAXBAS	Routing, length of weighting function		
<u><i>XAJ</i></u>			
WM	Areal soil moisture storage capacity		Lin et al. (2014)
B	The exponent of the soil moisture storage capacity curve		
KE	Ratio of potential evapotranspiration to pan evaporation		
IMP	Ratio of the impervious to the total area of the basin		
X	Proportion of soil moisture storage capacity of the upper layer to WM		
Y	Proportion of soil moisture storage capacity of the lower layer to WM		
C	Coefficient of deep evapotranspiration		
SM	Areal mean free water capacity of the surface soil layer		
EX	Exponent of the free water capacity curve		
KI	Coefficient of the free water storage to interflow		
KG	Coefficient of the free water storage to ground flow		
N	Number of reservoirs in the instantaneous unit hydrograph		
NK	Common storage coefficient in the instantaneous unit hydrograph		
CI	Recession constant of the lower interflow storage		
CG	Recession constant of the groundwater storage		

Table 4. Assumptions and descriptions of regionalization methods used in this study.

Method	Equation	Assumption and Description	Application examples
Spatial proximity	$D_{td} = \sqrt{(x_t - x_d)^2 + (y_t - y_d)^2}$	<p>Closer basins show similar hydrological characteristics.</p> <p>The donor catchments are determined by the distance <math>D_{td}</math>. <math>x, y</math> shows the location information, which uses the Universal Transverse Mercator (UTM) coordinate system.</p>	Merz and Blöschl (2004), Oudin et al. (2008), Yang et al. (2018, 2019)
Physical similarity	$SI_{td} = \sum_{i=1}^k \frac{ CD_{d,i} - CD_{t,i} }{\Delta CD_i}$	<p>Similar attributes show similarly in terms of hydrological processes.</p> <p>The donor catchments are decided by the similarity index <math>SI_{td}</math>. <math>CD</math> is the catchment descriptor, shown in Table 2 in this study.</p>	Burn and Boorman (1993), Poissant et al. (2017), Yang et al. (2018, 2019)
Regression	$MP_j = f_j(CD_i)$	<p>A well-behaved relationship exists between the observable <math>CDs</math> and model parameters (<math>MP</math>), and the <math>CDs</math> used in regression provide information relevant to hydrological behavior at ungauged sites.</p> <p>The relationship (linear regression function), which is built on gauged basins, will be transferred to ungauged catchments.</p>	Young (2006), Oudin et al. (2008), Merz et al. (2006), Yang et al. (2018, 2019)

$t$ : target catchment  
 $d$ : donor catchment  
 $i$ :  $i$ th catchment descriptors  
 $k$ : total number of catchment descriptors  
 $j$ :  $j$ th model parameter

$CD$ : catchment descriptor. The climate indices in  $CDs$  varied from the calibration to verification period, others are assumed as constant.



Table 5 The tested regionalization methods in this study

Regionalization methods	Abbreviation
Spatial proximity methods with parameter average option	SP-par
Spatial proximity methods with output average option	SP-out
Physical similarity methods with parameter average option	Phy-par
Physical similarity methods with output average option	Phy-out
Principal Component Regression method	PCR

Table 6. Average model performance in terms of Pbias, NSE and NSElog over the tested catchments in the split-sample test.

		calibration		verification	
		1980-1989	2006-2015	2006-2015	1980-1989
Pbias	GR4J	-0.81	-0.49	-4.37	-2.32
	WASMOD	2.61	3.15	-0.54	3.26
	HBV	-1.62	-1.49	-3.69	-3.90
	XAJ	-2.34	-1.69	-3.48	-1.80
NSE	GR4J	0.76	0.76	0.67	0.66
	WASMOD	0.77	0.76	0.68	0.67
	HBV	0.77	0.76	0.65	0.65
	XAJ	0.79	0.78	0.72	0.71
NSElog	GR4J	0.74	0.75	0.39	0.37
	WASMOD	0.67	0.71	0.58	0.55
	HBV	0.37	0.51	0.28	0.33
	XAJ	0.51	0.65	0.52	0.55

Table 7. The number of donor catchments providing the best performance for each regionalization method and hydrological model in the leave-one-out cross validation.

	GR4J	WASMOD	HBV	XAJ
SP-par	7	4	8	2
SP-out	10	9	8	9
Phy-par	3	2	2	3
Phy-out	3	5	5	3