

Regionalization in global hydrological models and its impact on runoff simulations: A case study using WaterGAP3 (v 1.0.0)

Jenny Kupzig¹, Nina Kupzig², Martina Flörke¹

¹Institute of Engineering Hydrology and Water Resources Management, Ruhr-University, 44801, Bochum, Germany

²Faculty of Management and Economics, Ruhr-University, 44780, Bochum, Germany

Correspondence to: Jenny Kupzig (jenny.kupzig@rub.de)

Abstract:

Valid simulation results from global hydrological models (GHMs), such as WaterGAP3, are essential to detecting hotspots or studying patterns in climate change impacts. However, the lack of worldwide monitoring data makes it challenging to adapt GHMs' parameters to enable such valid simulations globally. Therefore, regionalization is necessary to estimate parameters in ungauged basins. This study presents the results of regionalization methods for the first time applied on the GHM WaterGAP3. It aims to provide insights into (1) selecting a suitable regionalization method and (2) evaluating its impact on runoff simulation. In this study, four regionalization methods have been identified as appropriate for WaterGAP3. These methods span the full spectrum of methodologies, i.e., regression-based methods, physical similarity, and spatial proximity, using traditional and machine learning-based approaches. Moreover, the methods differ in the descriptors used to achieve optimal results, although all utilize climatic and physiographic descriptors. This demonstrates (1) that different methods use descriptor sets with varying efficiency and (2) that combining climatic and physiographic descriptors is optimal for regionalizing worldwide basins. Additionally, our research indicates that regionalization leads to spatially and temporally varying uncertainty in ungauged regions. For example, regionalization highly affects southern South America, e.g., leading to high uncertainties in the flood simulation of the Río Deseado. The local impact of regionalization propagates through the water system, also affecting global estimates, as evidenced by a spread of 1,500 km³ yr⁻¹ across an ensemble of five regionalization methods in simulated global runoff to the ocean. This discrepancy is even more pronounced when using a regionalization method deemed unsuitable for WaterGAP3, resulting in a spread of 4,208 km³ yr⁻¹. This significant increase highlights the importance of carefully choosing regionalization methods. Further research is needed to enhance the understanding of the methods' robustness on a global scale.

1. Introduction

Global hydrological models (GHMs) are developed and applied worldwide, e.g., to detect hotspots and examine patterns of climate change impacts on the terrestrial water cycle (e.g., Barbarossa et al., 2021; Boulange et al., 2021). Valid model results are a prerequisite to draw robust conclusions. For valid modeling results, it is beneficial to adjust the parameter values to adapt the models to different basin processes (Gupta et al., 1998). This adaptation is usually modified and evaluated (in a loop) by comparing the simulated model output, often discharge, with the monitored data. However, this parameter adjustment for GHMs is challenging due to the lack of global monitoring data. Consequently, parameter adjustment for GHMs can be based not only on monitored data (i.e., calibration) but also on estimating parameter values for ungauged basins (i.e., regionalization).

38 Regionalization defines the estimation of model parameters for ungauged basins (Oudin et al., 2008), usually based
39 on information from gauged basins (Oudin et al., 2010). Regionalization methods generally follow the same prin-
40 ciple: basin characteristics (e.g., physiographic and/or climatic) are linked to hydrological characteristics and can
41 thus be used to estimate parameter values. Various regionalization methods exist, and no overall preferred method
42 has been found (Ayzel et al., 2017; Pool et al., 2021). In contrast, the optimal regionalization method may differ,
43 for example, regarding available information (Pagliero et al., 2019) or model structures (Golian et al., 2021).
44 Therefore, different methods should be tested to find an optimal regionalization method for a specific use case
45 (e.g., Qi et al., 2020).

46 Evaluation is needed to assess different regionalization methods. The evaluation of regionalization methods is
47 particularly challenging because they are usually applied when there is a lack of monitoring data. Therefore, re-
48 gionalization studies often treat gauged basins as "ungauged" and perform leave-one-out cross-validation (e.g.,
49 Chaney et al., 2016) or split-sample tests (e.g., Beck et al., 2016; Nijssen et al., 2000; Yoshida et al., 2022). While
50 at the mesoscale, this evaluation is already an integral part (e.g., McIntyre et al., 2005; Parajka et al., 2005; Oudin
51 et al., 2008; Yang et al., 2020), this is sometimes not the case in global or continental studies (e.g., Müller Schmied
52 et al., 2021; Widén-Nilsson et al., 2007). Another reasonable evaluation strategy is the concept of benchmark-to-
53 beat (Schaeffli & Gupta, 2007; Seibert, 2001). Applying a benchmark-to-beat supports a comprehensive evaluation
54 of whether a new approach is functional, e.g., better than a straightforward and thus transparent method or better
55 than a predecessor. To the authors' knowledge, such a benchmark-to-beat has never been used to evaluate innova-
56 tions in regionalization at a global scale.

57 In general, regionalization methods can be divided into two categories based on the parameter estimation strategy:
58 (1) regression-based and (2) distance-based (He et al., 2011). Regression-based methods derive the relationship
59 between basin characteristics and model parameters through fitted regression models. These mathematically de-
60 fined relationships are further applied to estimate model parameters of ungauged basins (e.g., Kaspar, 2004; Müller
61 Schmied et al., 2021). A significant drawback of regression-based regionalization is the difficulty of incorporating
62 parameter interdependencies (Poissant et al., 2017), as regression-based approaches often assume that the depend-
63 ent variables, i.e., the model parameters, are not correlated (Wagener et al., 2004). Distance-based approaches
64 transfer complete parameter sets from similar or nearby donor basins to ungauged basins (e.g., Beck et al., 2016;
65 Nijssen et al., 2000; Widén-Nilsson et al., 2007). Using an ensemble of donor basins, e.g., by averaging the pa-
66 rameter values or model outputs, can improve the performance of such methods (e.g., Arsenault & Brissette, 2014).
67 A significant disadvantage of such methods is the clustering problem of ungauged basins, i.e., the unequal distri-
68 bution of gauging stations worldwide (Krabbenhoft et al., 2022). Thus, basins exist where distance-based ap-
69 proaches will use incomparable basins to transfer parameter values due to the lack of close basins.

70 Recent advances have implemented machine learning-based techniques in the context of regionalization. For ex-
71 ample, Chaney et al. (2016) used regression trees as an alternative to least squares regression to estimate parameter
72 values in ungauged basins. Pagliero et al. (2019) explored supervised and unsupervised clustering methods to
73 define the similarity of basins to transfer parameter sets. To the authors' knowledge, no study has compared several
74 traditional regionalization methods with machine learning-based methods for a GHM on a global scale.

75 Some regionalization methods do not make a clear distinction between calibration and regionalization. For exam-
76 ple, Arheimer et al. (2020) applied a basin grouping beforehand. Then, they jointly calibrated the group members
77 to define representative parameter sets. Subsequently, the representative parameter sets are transferred to other

78 basins based on grouping rules. Another approach defines so-called transfer functions (Samaniego et al., 2010)
79 and calibrates meta-parameters instead of the model parameter values (Beck et al., 2020; Feigl et al., 2022). These
80 methods, where regionalization is part of the calibration process, often require a change in the calibration process
81 itself, which is challenging for GHMs (Schweppe et al., 2022), for example, due to a lack of code flexibility (e.g.,
82 Cuntz et al., 2016).

83 This study proposes an improved regionalization method for the state-of-the-art GHM WaterGAP3 (Eisner, 2016).
84 It compares traditional regionalization methods with machine learning-based methods and uses a benchmark-to-
85 beat and an ensemble of split-sample tests to evaluate the applied methods. Further, global runoff simulations are
86 compared to analyze the impact of regionalization methods. The overall research topic is evaluating and selecting
87 regionalization methods for a GHM. Specifically, the study has two objectives. It aims

- 88 (1) to propose an improved regionalization method for WaterGAP3 and
- 89 (2) to evaluate the impact of regionalization methods on global runoff simulations.

90 **2. Data and Methods**

91 **2.1 The Model: WaterGAP3**

92 The GHM WaterGAP3 simulates the terrestrial water cycle, including the main water storage components and a
93 simple storage-based routing algorithm. It is a fully distributed model that operates on a five arcmin grid and
94 simulates at a daily time step. A more detailed description of the model can be found in Eisner (2016).

95 In WaterGAP3, most model parameter values are set a priori, e.g., using look-up tables for albedo or rooting depth.
96 Only one parameter, γ , is calibrated, which is part of the soil moisture storage in which runoff generation processes
97 are present. The model equation for γ , which originates from the HBV-96 model (Lindström et al., 1997), is given
98 in Eq. (1). Generally, higher values of γ lead to lower runoff volumes, while lower values of γ lead to higher runoff
99 volumes. The model parameter is calibrated per basin within the range of 0.1 and 5. The objective function of the
100 calibration is to minimize the deviation between the mean annual simulated and observed river discharge, i.e., the
101 calibration aims to reduce the error in discharge volume. Given the monotonic relationship between the model's
102 parameter and the optimization function, a simple search algorithm is applied: The parameter space is divided into
103 rectangles, which are subsequently subdivided into smaller rectangles depending on the direction γ should be
104 modified to achieve closer alignment with the optimization target. The calibration results in one calibrated γ value
105 between 0.1 and 5 per basin. After the calibration, a correction is applied to account for high errors in the mass
106 balance, e.g., due to inaccuracies in global meteorological forcing products. This correction is only applicable on
107 gauged basins. It is, therefore, neglected in this study.

$$108 \quad R = P_t \cdot \left(\frac{S_s}{S_{s,max}} \right)^\gamma \quad (1)$$

109 where R is the daily runoff, P_t is the daily throughfall, S_s is the actual soil storage, $S_{s,max}$ is the maximal soil
110 storage (given as a global map in Appendix A), and γ is the calibration parameter.

111 Traditionally, the regionalization process in WaterGAP3 is a simple multiple linear regression (MLR) approach to
112 estimate the calibration parameter γ for ungauged basins (e.g., Döll et al., 2003; Kaspar, 2004). The drawback of
113 MLR regarding parameter interaction can be neglected: As there is only one parameter to estimate, parameter

114 interference does not exist. Instead, the approach offers the advantage of a lightweight, transparent application that
115 can be quickly revised and adapted.

116 **2.2 Model Data**

117 WaterGAP3 requires various input data, such as soil information, topography, or information on open freshwater
118 bodies. This study uses the same input data as Kupzig et al. (2023). For meteorological forcing, we use the global
119 data set EWEMBI (Lange, 2019). This data product includes daily global forcing data with a spatial resolution of
120 0.5 degrees (latitude and longitude) that covers a period from 1979 to 2016. Specifically, WaterGAP3 uses the
121 following forcing information from the EWEMBI data set as input:

- 122 • daily mean temperature,
- 123 • daily precipitation,
- 124 • daily shortwave downward radiation, and
- 125 • daily longwave downward radiation.

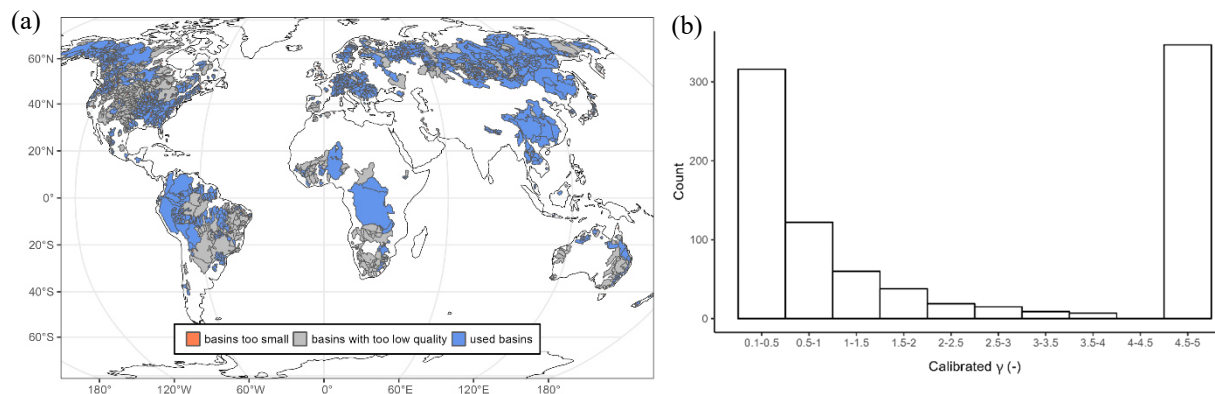
126 The WaterGAP3 calibration requires observed monthly river discharge data. This discharge data is subsequently
127 transformed into annual discharge sums and used as a benchmark in the calibration procedure. In this study, we
128 used discharge data from 1,861 stations that were manually verified (Eisner, 2016). To get the best data available,
129 we have updated all available station data with recent data from The Global Runoff Data Center (GRDC, 2020).
130 All stations have at least five years of complete (monthly) station data between 1979 and 2016. For each station,
131 a contribution area, i.e., a basin, is defined with the gridded flow-direction information obtained from WaterGAP3,
132 based on the HydroSHEDS database (Lehner et al., 2008).

133 The 1,861 basins are calibrated using the above-described standard calibration approach for WaterGAP3. Follow-
134 ing the standard calibration procedure, some basins still have an insufficient model performance. In this context,
135 we define a monthly Kling-Gupta-Efficiency (KGE) below 0.4 or more than 20 % bias in monthly flow as insuf-
136 ficient model performance. We underscore the importance of minimizing the error in discharge volume by defining
137 it as an additional criterion corresponding to the optimization target during calibration. Basins not fulfilling the
138 defined conditions regarding bias and KGE are neglected in further analysis to avoid high parameter uncertainty
139 due to errors in input data, model structure, or discharge data affecting the analysis. Further, we have excluded all
140 basins with less than 5000 km² (inter-) basin size from the next upstream basin. We assume that this inter-basin
141 size is large enough to assume a certain degree of interdependency between nested basins. In total, 933 out of
142 1,861 basins are selected for regionalization (626 are neglected due to insufficient model performance, and 302
143 are neglected due to inadequate basin size).

144 Figure 1a depicts the worldwide calibrated basins, highlighting gauged and ungauged regions. Whereas most parts
145 of North and South America are gauged, Africa and Australia remain largely ungauged. A cluster of gauged basins
146 is in Central Europe and in Eastern Asia. Gauged regions with insufficient model performance are mainly in the
147 Mississippi River basin, Southern Africa, Australia, and large parts of Brazil. These regions are known to be chal-
148 lenging for GHMs (e.g., cf. Fig. 8b in Stacke & Hagemann, 2021).

149 Figure 1b shows the calibrated values for γ . It emerges that the calibrated values tend to be at the upper and lower
150 bounds of the parameter space. This behavior is already known (cf. Fig. 4b in Müller Schmied et al., 2021). A
151 brief sensitivity analysis and discussion of the calibration parameter are included in Appendix B. The results of

152 this analysis indicate that the clustering of the calibrated parameter value is not related to an inappropriate selection
 153 of the parameter bounds but instead to the absence or an insufficient representation of processes. Thus, the clus-
 154 tering of the calibrated values does not indicate an inadequate selection of the parameter bounds but highlights the
 155 necessity to improve the model structure and the calibration strategy for WaterGAP3. However, this study focuses
 156 solely on analyzing and implementing regionalization methods. It does not aim to enhance the model structure or
 157 to change the calibration procedure of WaterGAP3. Future studies are needed to achieve the latter, as WaterGAP3
 158 contains many hard-coded parameters or parameters defined by look-up tables that need to be analyzed to identify
 159 and adjust sensitive parameters more accurately during calibration. Initial steps in this direction have already been
 160 taken for WaterGAP2 in the form of a multivariate and multi-objective case study in the Mississippi River basin
 161 (Döll et al., 2024).



162 **Figure 1: (a) Map of calibrated basins, highlighting basins not used for regionalization due to insufficient model per-**
 163 **formance or inadequate basin size and (b) the histogram of the calibrated model parameter values of all used basins**
 164 **showing a cluster of parameter values at the parameter bounds.**

165 2.3 Basin Descriptors

166 This study uses basin descriptors as predictors to drive regression-based or distance-based regionalization ap-
 167 proaches. These basin descriptors are based on data used within the model simulation (as they are globally avail-
 168 able). They are aggregated to basin values using a simple mean method to have the same spatial resolution as the
 169 calibrated model parameter. Thus, in the case of nested basins, the inter-basin area is used to define the basin
 170 descriptors. The selection of the predictors, i.e., basin descriptors that support the estimation of γ , is crucial for
 171 regionalization methods (Arsenault & Brissette, 2014). Typically, this selection aims to obtain the most infor-
 172 mation with the least number of predictors to (1) improve the model quality and (2) limit over-parametrization. In
 173 this study, we use 12 basin descriptors to develop regionalization methods; nine of these descriptors are physio-
 174 graphic, while the remaining three are climatic (see Table 1). Most descriptors are not correlated (see Appendix
 175 C), i.e., we minimize redundant information (Wagener et al., 2004).

176 A descriptor subset is selected based on correlation analysis between basin descriptors and calibrated γ value and
 177 entropy assessment. Pearson's correlation coefficient detects linear correlation, and Spearman's Rho and Kendall's
 178 Tau detect a non-linear correlation. Shannon entropy (Shannon, 1948) measures the information gain of the pre-
 179 dictors explaining the calibrated γ value. The higher the information gain, the more valuable the basin descriptor
 180 is for explaining the variation in the calibrated γ value. The analysis directly evaluates the relationship between
 181 the calibrated parameter and the basin descriptors, as WaterGAP3 uses only one calibration parameter with a clear
 182 global optimum within the parameter space. An alternative would be to use flow characteristics to define the basis

183 for regionalization (e.g., Pagliero et al., 2019). We decided to use the calibrated parameter instead of flow charac-
 184 teristics as it does not need any further assumption on which flow characteristics determine the model's parameter.
 185 Statistical information of the evaluated basin descriptors and the corresponding correlation coefficients and infor-
 186 mation gain are listed in Table 1. The basin descriptors demonstrate a considerable degree of variability, e.g., the
 187 basin size ranges from 5000 km² to 3,112,480 km² with a median of 13,796 km². The mean temperature varies
 188 from -19 °C to 29 °C, and the sum of precipitation ranges from 213 mm to 5,716 mm. Although there is a high
 189 degree of variability in the analyzed basin descriptors, the basin descriptors exhibit low correlation coefficients
 190 with the calibrated values. For example, the permafrost coverage shows the strongest Pearson correlation of -0.37
 191 (and -0.50 for Spearman's Rho). The information gain indicates the same results as the correlation analysis, i.e.,
 192 the information gain is generally relatively low, and descriptors with a higher correlation tend to have a higher
 193 information gain. For example, the mean temperature exhibits the maximal information gain of 17.6 % and has
 194 the second-highest correlation coefficient with a Pearson correlation of 0.34.

195 **Table 1: Basin descriptors: statistical information, correlation, and entropy assessment. Selected physiographic and**
 196 **climatic basin descriptors are written in bold.**

	Basin Descriptor	Attribute Information				Entropy & Correlation			
		Min	Max	Mean	Median	IG (%) ¹	Pearson	Spearman	Kendall
physiographic	Soil Storage (mm)	12.405	610.469	220.805	195.778	13.07	-0.21	-0.15	-0.11
	Open Water Bodies (%)	0.000	63.960	5.521	1.812	5.65	-0.01	-0.08	-0.05
	Wetlands (%)	0.000	63.466	4.164	0.547	5.01	-0.02	-0.13	-0.09
	Size (km ²)	5000	3,112,480	37,572	13,796	1.42	-0.04	-0.04	-0.03
	Slope Class (-)	10.057	67.756	38.668	38.364	16.60	-0.31	-0.37	-0.27
	Altitude (m.a.s.l.)	30.239	4765.166	591.024	394.870	9.30	-0.18	-0.28	-0.20
	Sealed Area (%)	0.000	12.3	0.6	0.1	4.49	0.22	0.38	0.29
	Forest (%)	0.000	100.000	35.340	24.002	13.82	-0.25	-0.18	-0.14
	Permafrost & Glacier (%)	0.000	95.000	16.662	0.000	13.12	-0.37	-0.50	-0.40
climate	Mean Temperature(°C)	-18.848	28.823	7.720	7.707	17.56	0.34	0.41	0.30
	Yearly Precipitation (mm)	213.6	5,716.3	996.5	779.5	9.23	0.02	0.21	0.14
	Yearly Shortwave Down-ward Radiation (Wm⁻²)	1,050.6	3,043.2	1,857.9	1,759.7	15.79	0.31	0.33	0.24

¹Information gain is given in percentage of total information content in γ after Shannon (1948)

197 In contrast to the findings of Wagener and Wheater (2006), the correlation coefficients between the basin de-
 198 scriptors and the calibrated values are relatively low, indicating a weak relationship. One potential explanation for
 199 this discrepancy is that Wagener and Wheater (2006) used a smaller number of basins in southeast England, with
 200 limited versatility (e.g., regarding climate and seasonality) compared to the 933 worldwide basins used in this
 201 study. Studies using a large number of basins likely tend to find a lower correlation between catchment attributes
 202 and model parameters (Merz et al., 2004). Moreover, the clustered calibrated γ values at the bounds of the valid
 203 parameter space may disturb the results of this analysis. As the calibrated value masks the effect of multiple sources
 204 of errors, such as uncertainty in the input data, model structure, or varying hydrological processes, finding a mean-
 205 ingful relationship between catchment characteristics and calibrated values is challenging.

206 Because the basis for the descriptor selection seems uncertain, given the low correlation and the named constraints,
 207 we additionally run the regionalization methods with all descriptors to evaluate the descriptor selection. Further
 208 on, to ascertain the advantage of integrating climatic descriptors, we run the regionalization methods using either
 209 physiographic or climatic descriptors. In total, we used four groups of basin descriptors to implement the region-
 210 alization methods:

- 211 • "cl": all three climatic descriptors,
- 212 • "p": all nine physiographic descriptors,
- 213 • "p+cl": all 12 descriptors, and
- 214 • "subset": two correlated climatic descriptors (mean temperature, annual shortwave radiation) & three
- 215 correlated physiographic descriptors (slope class, forest %, permafrost %).

217 2.4 Regionalization Methods

218 In our study, we test several traditional and machine learning-based regionalization methods against each other
 219 and a defined benchmark-to-beat to find suitable regionalization methods for WaterGAP3. At the global scale,
 220 regionalization is particularly challenging due to (1) the lack of high-quality data, (2) the diversity of dominant
 221 hydrological processes in basins, and (3) the high computational demands of the models. Therefore, a robust re-
 222 gionalization method that applies to a wide variety of basins and is not computationally demanding should be
 223 selected for a global application.

224 We test three common traditional approaches and two machine learning-based approaches using the concepts of
 225 spatial proximity, physical similarity, and regression-based methods. As WaterGAP3's model calibration is very
 226 rigid and has only one parameter, it is not feasible to implement and test regionalization methods that incorporate
 227 regionalization into the calibration process, such as transfer functions. In addition, we avoid high computational
 228 demands as all evaluated methods are applicable after the calibration, i.e., without running the model.

229 As the calibration of WaterGAP3 results in a parameter distribution with a cluster of parameter values at the
 230 parameter bounds, we implement a so-called "tuning" to introduce information about the parameter space into
 231 regionalization. In detail, we apply a simple threshold-based approach to shift the regionalized parameter values
 232 to the extremes, i.e., $\gamma_{est} < \gamma_1 \rightarrow \gamma_{reg} = 0.1$ and $\gamma_{est} > \gamma_2 \rightarrow \gamma_{reg} = 5.0$. The thresholds γ_1 and γ_2 are defined
 233 by applying the k-means algorithm with three centers to the calibrated parameter values. This clustering results in
 234 three clusters: one for low, one for medium, and one for high γ values. Subsequently, γ_1 refers to the highest γ
 235 value of the low cluster and γ_2 refers to the lowest γ value of a high cluster.

236 To evaluate the regionalization methods, we implement an ensemble of split-sample tests. Specifically, we ran-
 237 domly split the basins into 50 % gauged (for training) and 50 % pseudo-ungauged (for testing). The split has a
 238 relatively high percentage of pseudo-ungauged basins, accounting for many missing gauges worldwide. We fit the
 239 methods and apply them to the training and testing data sets. The split-sample test is repeated 100 times by ran-
 240 domly splitting the basins to account for sampling effects.

241 As there is only one calibration parameter, γ , this parameter has a global optimum per basin. Consequently, the
 242 quality of training and testing is directly assessed by the deviation between the regionalized and the calibrated
 243 value for γ . The closer the regionalized values are to the calibrated ones, the more accurate the prediction. We
 244 assess the prediction accuracy by the logarithmic version of the mean absolute error (logMAE) to account for the
 245 decreasing sensitivity of γ for higher values (see Appendix B). The lower the logMAE, the better the prediction; a
 246 zero value in logMAE expresses no error. The regionalization method is robust if the prediction accuracy is similar
 247 in training and testing. A generally good performance, i.e., small logMAE values, indicates that the regionalization
 248 method suits WaterGAP3. The comparison of γ values enables applying a wide range of regionalization methods
 249 and sets of descriptors, as no computationally intensive model simulation is required. However, it assumes that
 250 deviations in γ lead, in turn, to deviations in discharge, which is only partially true because of varying parameter

251 sensitivity in basins (e.g., Kupzig et al., 2023). To validate that the logMAE is a sufficient approximator for the
252 regionalization performance in WaterGAP3, we use one representative split-sample from the ensemble to compare
253 the accuracies in simulated discharge for different regionalization methods.

254 **Regression-based methods**

255 The traditionally used regionalization approach of WaterGAP3 is a regression-based MLR. As the benchmark-to-
256 beat, we use the regionalization approach from WaterGAP2.2d defined in Müller Schmied et al. (2021). We con-
257 sider it a suitable benchmark-to-beat given that WaterGAP2 has a model structure and calibration process that is
258 very similar to WaterGAP3. The main difference between these models is that WaterGAP2 simulates at 0.5° spatial
259 resolution. The benchmark-to-beat consists of "a multiple linear regression approach that relates the natural loga-
260 rithm of γ to basin descriptors (mean annual temperature, mean available soil water capacity, fraction of local and
261 global lakes and wetlands, mean basin land surface slope, fraction of permanent snow and ice, aquifer-related
262 groundwater recharge factor)". (Müller Schmied et al., 2021) We fit this regression model to our data and define
263 the quality of this approach as the benchmark-to-beat. Moreover, we test an independent MLR approach without
264 using the logarithmical scaling of γ and using the above-defined sets of basin descriptors. For MLR and the bench-
265 mark-to-beat, we use the `lm()` function of the R package `stats` (R Core Team, 2020). After applying the regression
266 model, we adjust the estimated parameter values to ensure that the estimated values range between 0.1 and 5.

267 Furthermore, a machine learning-based method, random forest (RF), is tested for regionalization as an alternative
268 to MLR. Here, we implement the random forest algorithm with the `randomForest()` function from the R package
269 `randomForest` (Liam & Wiener, 2002), which is based on Breimann (2001). The algorithm uses an ensemble of
270 decision trees, making the decision human-like. It is relatively robust because it incorporates random effects into
271 the training process. To implement this randomness, we define the algorithm as one that can choose between two
272 randomly selected predictors at each node, using an ensemble of 200 trees.

273 **Physical Similarity**

274 As the traditional physical similarity approach, we use Similarity Indices (in the following named with SI), apply-
275 ing the methodology proposed by Beck et al. (2016). The SI (see Eq. (2)) are derived using the defined basin
276 descriptors sets, and the parameter of the most similar basin is transferred to the pseudo-ungauged basin. Addi-
277 tionally, we use an ensemble of basins to control whether an ensemble-based approach leads to more robust results.
278 The optimal number of donor basins may vary between research regions and hydrological models (Guo et al.,
279 2020). Here, we use ten donor catchments (noted with "ensemble") based on Beck et al. (2016) and McIntyre et
280 al. (2005). Further, we apply a simple mean method for the ensemble-based prediction to aggregate the ensemble
281 of γ values into one predicted parameter value.

$$282 \quad S_{i,j} = \sum_{p=1}^n \frac{|Z_{p,i} - Z_{p,j}|}{IQR_p} \quad (2)$$

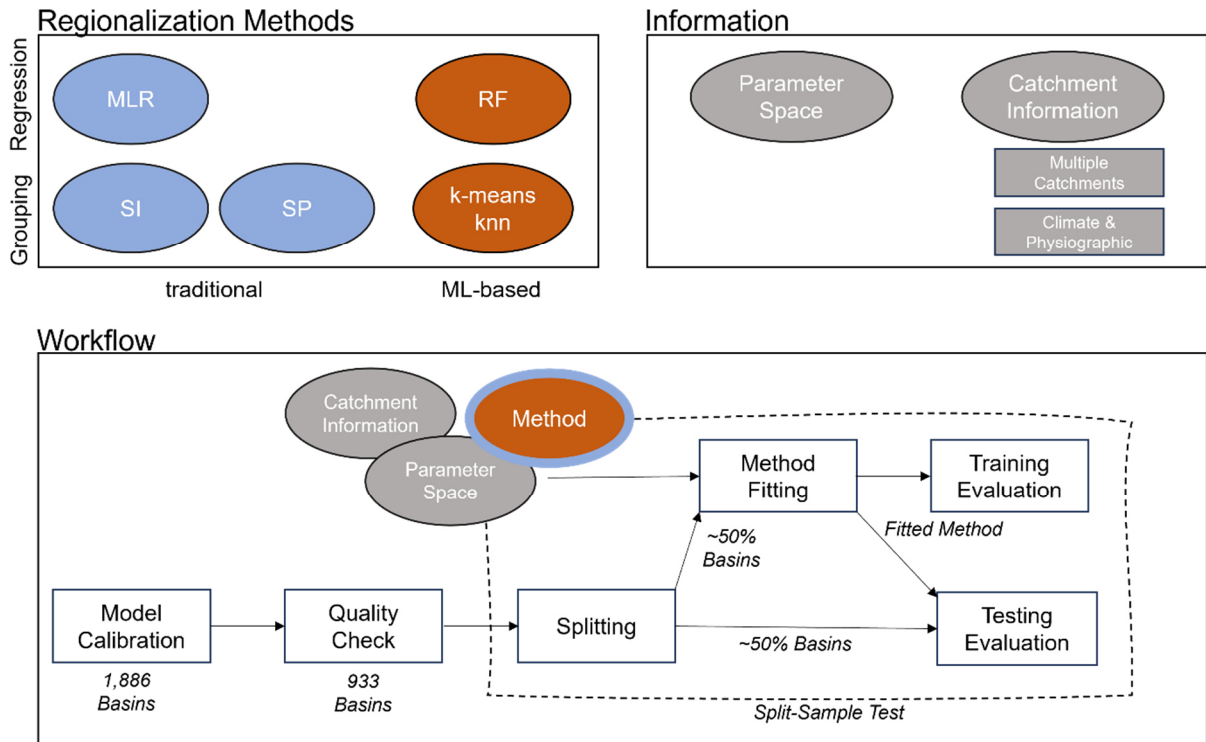
283 where $S_{i,j}$ is the Similarity Index between basin i and basin j , $Z_{p,j}$ is the basin descriptor p for basin j , IQR_p is the
284 interquartile range for basin descriptor p among all (gauged) basins, and n is the number of all basin descriptors
285 used.

286 As an alternative machine learning-based approach, we apply a simple k-means algorithm. We selected the k-
287 means algorithm because it is one of the most widely used clustering algorithms (Tongal & Sivakumar, 2017). It

288 is easy to understand and use. The algorithm `kmeans()` is implemented in the R base package `stats`. It aims to
 289 maximize variation between groups and minimize variation within groups. The number of clusters to use is deter-
 290 mined by multiple indices calculated with the R package `NbClust` (Charrad et al., 2014). For all 933 basins and
 291 the defined sets of basin descriptors, most indices defined three as the optimal number of clusters. Accordingly,
 292 we use three clusters to generate the groups of basins. As different scales of the predictor values can affect the
 293 clustering, a rescaling with min-max-normalization (see Eq. (3)) is performed on the training set and applied to
 294 the testing set. After the grouping, the mean γ value is assigned as a representative calibrated value to the corre-
 295 sponding basin group. To estimate the corresponding group for a pseudo-ungauged basin, the `knn` algorithm is
 296 used, and the representative γ value of the group is assigned to the pseudo-ungauged basin. This algorithm is
 297 implemented by the `knn()` function of the R package class (Venables & Ripley, 2002). Since the k-means method
 298 is less flexible than SI, we implement a highly flexible version, using the `knn` algorithm directly to define the donor
 299 basin most similar to each ungauged basin. Using the `knn` algorithm directly, we test how beneficial it is to create
 300 groups of similar basins using the `kmeans` algorithm and regionalize the parameter with a representative mean
 301 value.

$$302 \quad Z'_{p,j} = \frac{Z_{p,j} - \min_{j \rightarrow m}(Z_{p,j})}{\max_{j \rightarrow m}(Z_{p,j}) - \min_{j \rightarrow m}(Z_{p,j})} \quad (3)$$

303 where $Z'_{p,j}$ is the normalized basin descriptor p for basin j , $Z_{p,j}$ is the basin descriptor p for the basin j , m is the
 304 number of (gauged) basins.



305
 306 **Figure 2: Experimental setup of the study: regionalization methods, used modifications and information, and the gen-**
 307 **eral workflow (MLR: Multiple Linear Regression, SI: Similarity Indices, SP: Spatial Proximity, RF: RandomForest).**

308 **Spatial Proximity**

309 The spatial proximity approach is one of the easiest to regionalize parameter values. However, it is also often
310 criticized that nearby basins do not necessarily have the same hydrological behavior (Wagener et al., 2004). Fur-
311 thermore, its performance depends on the density of the network of gauged basins (Lebecherel et al., 2016). The
312 dependency on network density is particularly challenging for global applications where large parts of the world
313 are ungauged (e.g., northern Africa). Nevertheless, the approach has been successfully applied in other studies
314 (e.g., Oudin et al., 2008; Qi et al., 2020), even globally (Widén-Nilsson et al., 2007). Here, we take the distance
315 between the centroids of the basins as the reference for the spatial distance between basins, as done by others
316 (Oudin et al., 2008). We use the abbreviation SP in the text below to refer to the spatial proximity approach.
317 Figure 2 provides an overview of the applied regionalization methods and information used for the experimental
318 setup.

319 **3. Results and Discussion**

320 **3.1 Evaluating the effect of tuning**

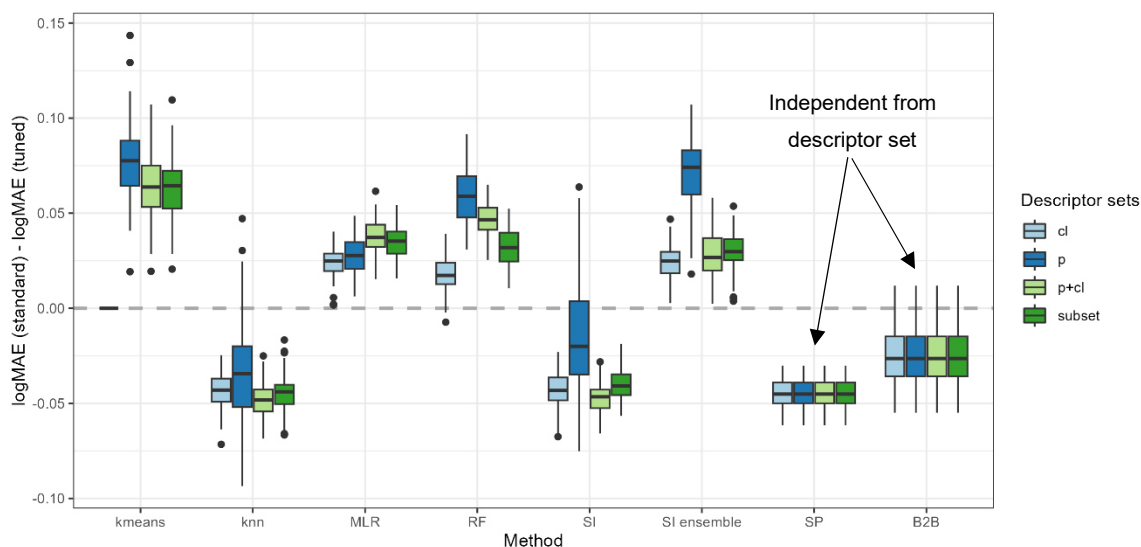
321 First, the impact of the tuning approach on the regionalization approaches is evaluated. Therefore, Fig. 3 depicts
322 the differences in logMAE between the standard and tuned approaches in testing, i.e., using the pseudo-ungauged
323 basins. A positive difference in logMAE indicates an increase in accuracy, whereas a negative difference indicates
324 a decrease in accuracy due to the tuning.

325 Using the tuning thresholds of about 1.1 and 3.4 for γ_1 and γ_2 , respectively, enhances the predictive accuracy for
326 kmeans, MLR, RF, and the ensemble approach of SI. The most remarkable improvement for kmeans, RF, and SI
327 ensemble is achieved when all physiographic descriptors are used as input (mean improvement of 0.077, 0.058,
328 and 0.071, respectively). MLR shows the most significant improvement when using all available descriptors (mean
329 improvement of 0.038). In contrast, the tuning decreases the performance for knn, SI, and SP, with a mean degra-
330 dation between -0.02 and -0.05. Unlike the enhanced regionalization techniques, these methods transfer single-
331 basin information to ungauged regions. Thus, the tuning disturbs the use of single-basin information yet simulta-
332 neously enhances the performance of methods that transfer multi-basin information. The disturbance or improve-
333 ment is probably related to the capability of the methods representing the clustering of parameter values at the
334 extremes: Whereas the multi-basin information transfer implies a smoothing and thus suffers from a lack of rep-
335 resenting the extremes, the single-basin information transfer exhibits no such a smoothing.

336 The exception from the above-defined rule is the benchmark-to-beat approach. The benchmark-to-beat is the only
337 approach that uses logarithmic scaled γ values when fitting the model. This logarithmic transformation leads to an
338 increase in estimating small values. Thus, when the benchmark-to-beat is tuned, more basins with higher calibrated
339 γ values receive low estimates. The tuning intensifies this effect, leading to a decrease in the accuracy of the
340 logMAE from the standard to the tuned version. Thus, for models using logarithmical transformed γ values, the
341 defined thresholds for the tuning are not appropriate.

342 Applying knowledge of the optimal parameter space enhances the quality of regionalization for methods transfer-
343 ring multi-basin information in case the tuning thresholds are appropriate. This positive effect is not surprising, as
344 incorporating a priori information about parameter distribution strengthens parameter estimation (e.g., described
345 in Tang et al. (2016) using the Bayes Theorem). However, for single-basin transfer, which already represents the

346 parameter space well, i.e., the clustering of γ at the extremes, the tuning disturbs the performance. This indicates
 347 that such tuning needs to be cautiously introduced as there is the risk of decreasing the accuracy of regionalization.



348
 349 **Figure 3: Changes in performance between standard and tuned versions for all applied regionalization approaches.**
 350 **Positive values indicate an improvement related to the tuning.**

351 3.2 Evaluating descriptor subsets & algorithm selection

352 Different descriptor sets yield different performances in regionalizing γ . Table 2 shows the median of all logMAE
 353 values for the testing. For a complete overview of the results of the split-sample test ensemble, see Appendix D.
 354 Evaluating Table 2 reveals that the selected subset or all descriptors consistently yield the best performance across
 355 all regionalization methods. In both variants of the ensemble approach of SI, the tuned version of the no-ensemble
 356 approach of SI, and the standard version of RF, the selected subset yields the best results. For all other methods,
 357 using all descriptors yields the best results. Hence, all methods perform best when combining climatic and physi-
 358 ographic descriptors. This benefit of using climatic and physiographic descriptors is consistent with others that
 359 often apply a combination of climatic and physiographic descriptors, achieving optimal regionalization results
 360 (e.g., Oudin et al., 2008; Reichl et al., 2009).

361 The machine learning-based approaches seem to benefit most when using more information displaying an im-
 362 provement for all methods (knn, kmeans, and RF) and both variants (standard and tuned) ranging from "cl", "p",
 363 "subset" to "p+cl". This is not surprising as machine learning is developed to deal with big data sets. The traditional
 364 methods MLR and SI do not exhibit such a distinct pattern. The (weakly) correlated subset of climatic and physi-
 365 ographic descriptors yields the best results for SI. As utilizing all descriptors decreases the performance slightly,
 366 the results indicate that uncorrelated descriptors may disturb the performance of this approach. For MLR, the
 367 meaning of physiographic information is highest, resulting in the best ("p+cl") and second best ("p") results. The
 368 disparate performance of the regionalization methods when using different descriptor sets indicates that different
 369 methods use descriptor sets with varying efficiency. It also emphasizes that the selection of descriptors impacts
 370 the regionalization method's results, as noted by others (Arsenault & Brissette, 2014). Consequently, the above-
 371 performed analysis defining a descriptor subset lacks universal validity as methods exist where the defined subset
 372 is outperformed. Instead, the validity of this approach is most closely aligned with the SI approaches.

373 Although the algorithms kmeans and knn are similar, they yield considerably different performances in Table 2.
 374 As knn shows a logMAE of 0.432 at best, the kmeans algorithm performs poorly, resulting in the best logMAE of
 375 0.472. This indicates that applying the kmeans clustering algorithm to transfer averaged parameters is inappropriate
 376 for WaterGAP3. This may be attributed to the reduced flexibility of the approach, which entails estimating
 377 only three γ values due to the optimal, though limited, number of centers. The ensemble SI approach consistently
 378 outperforms the no-ensemble SI approach in almost all variants. The positive effect of an ensemble approach for
 379 SI has already been noted (Oudin et al., 2008). Therefore, it is recommended that the number of donor basins
 380 derived from the literature be adopted in future applications to be optimal for WaterGAP3, likely resulting in
 381 higher performance.

382 Only a few regionalization methods outperform the benchmark-to-beat. The best descriptor sets of tuned MLR,
 383 RF, and SI ensemble approach have a logMAE of 0.427, 0.403, and 0.409, respectively. The standard version of
 384 knn ("p+cl") and SP yield 0.432 and 0.454 in logMAE, respectively. Additionally, two variants of the standard SI
 385 approaches outperform the benchmark-to-beat yet exhibit inferior results compared to the selected tuned approach.
 386 All other regionalization methods show higher logMAE values than the benchmark-to-beat. These methods are
 387 considered insufficient in terms of performance to regionalize γ in WaterGAP3. As the benchmark-to-beat outper-
 388 forms all kmeans approach variants, it is deemed unsuitable for regionalizing γ for WaterGAP3 and, therefore,
 389 excluded from further analysis.

390 **Table 2: Median logMAE of 100 split-samples for pseudo-ungauged basins, i.e., in testing, for all regionalization meth-**
 391 **ods applying four sets of descriptors for a) the standard version and b) the tuned version. The bold numbers indicate a**
 392 **better performance than the benchmark-to-beat. Thicker edges mark best-performing variants, which are chosen for**
 393 **further analysis. Grey-shaded cells indicate worst-performing variants, which were taken to validate the assumption**
 394 **that lower logMAE values result in lower KGE values.**

(a)

test (median)	MLR	RF	SI		kmeans	knn	SP	B2B
			no ens.	ensemble				
cl	0.552	0.483	0.496	0.483	0.619	0.501	0.454	0.461
p	0.479	0.465	0.487	0.480	0.551	0.477		
p+cl	0.464	0.464	0.454	0.462	0.534	0.432		
subset	0.488	0.488	0.461	0.439	0.539	0.467		

(b)

test* (median)	MLR	RF	SI		kmeans	knn	SP	B2B
			no ens.	ensemble				
cl	0.529	0.467	0.537	0.459	0.619	0.546	0.502	0.488
p	0.441	0.416	0.532	0.455	0.515	0.521		
p+cl	0.427	0.403	0.503	0.435	0.472	0.480		
subset	0.453	0.408	0.501	0.409	0.477	0.509		

395 The well-performing SP on a global scale is surprising as the distances between basins are potentially long, and
 396 hydrological processes may strongly vary. It is probably beneficial for the SP approach that γ comprises all kinds
 397 of errors, e.g., spatially localized errors in global forcing products (e.g., Beck et al., 2017 reported errors for arid
 398 regions in the precipitation product) or inaccurately represented processes for larger regions. Thus, the estimation
 399 of γ might be appropriate, but not because of the same hydrological behavior but due to the same kind of errors.

400 The RF approach is outstanding, as it shows a massive loss in performance from training to testing (see Appendix
401 D). In detail, the logMAE in testing is about twice the logMAE in training. In comparison, other methods show
402 results from 95.6 % to 101.4 %. This performance loss indicates that RF is not a robust regionalization method for
403 WaterGAP3. Other studies that reported the good performance of RF for regionalization have not investigated the
404 stability of the performance from training to testing (Golian et al., 2021; Wu et al., 2023). Likely, the mathematical
405 problem of predicting the calibrated parameter for WaterGAP3, with all its challenges (e.g., tailored parameter
406 space, clustered calibrated parameter, and incorporation of many sources of errors), cannot be adequately solved
407 by RF. Thus, although RF is known to be especially robust among other machine learning-based techniques, it
408 shows symptoms of over-parameterization. This indicates that the algorithm is too flexible and adjusts to noise in
409 the data, missing the underlying systematic. This lack of robustness is particularly disadvantageous since, for Wa-
410 terGAP3, regionalization is applied globally, requiring regionalizing large parts of the world. In consequence, the
411 RF approach is left out from further analysis and defined as not suitable to regionalize γ for WaterGAP3.

412 3.3 Performance of selected algorithm in pseudo-ungauged basins

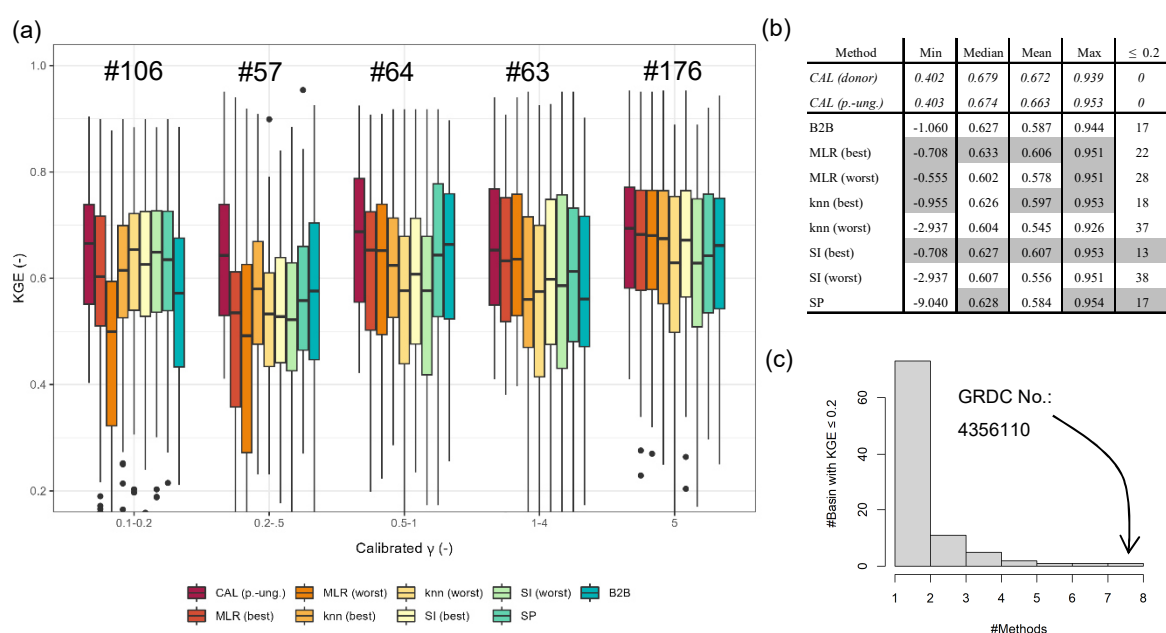
413 To avoid the high risk of sampling effect when applying the split-sample test, we conduct an ensemble of 100
414 split-sample tests analyzing the median of logMAE between regionalized and calibrated values as an indicator for
415 performance. Directly using the differences in regionalized and calibrated values is only meaningful when the
416 calibrated value represents the global optimum. As this is often not the case, e.g., due to equifinality, the perfor-
417 mance of regionalization methods is usually assessed by the accuracy of simulated discharge (e.g., Samaniego et
418 al., 2010; Arsenault & Brissette, 2014). Because WaterGAP3 requires computationally intensive simulations, run-
419 ning WaterGAP3 for all 100 split-sample tests for the selected methods is not feasible. Therefore, we select a
420 single representative split-sample to assess the quality of representing the discharge in the pseudo-ungauged basins
421 using regionalized γ values. The representative split-sample leads to comparable logMAE values to the corre-
422 sponding median of the ensemble for all regionalization methods. For the evaluation, WaterGAP3 was run for the
423 same period used in calibration (from 1979 to 2016), with the first year simulated ten times to allow for model
424 warm-up. Using this period ensures the availability of sufficient data for the evaluation (see Chapter 2.2). Further-
425 more, the differences between the monthly simulated and observed discharge are assessed using the KGE.

426 To evaluate the KGE, we select the best-performing methods that outperform the benchmark-to-beat: tuned MLR
427 "p+cl", knn "p+cl", tuned SI ensemble "subset", and SP (see Table 2). For the sake of simplicity, we further mark
428 them with "(best)". Additionally, we select three poorly performing variants to validate the assumption that meth-
429 ods resulting in higher logMAE values tend to result in lower KGE values, i.e., lower accuracy of simulated dis-
430 charge. These methods are tuned SI "cl" (logMAE: 0.537), tuned knn "cl" (logMAE: 0.546), and MLR "cl" (log-
431 MAE: 0.552). Further, we denote these methods with "worst". Applying the selected methods and the benchmark-
432 to-beat method results in eight estimates of γ for the pseudo-ungauged basins, whose performance is further eval-
433 uated in terms of simulated discharge accuracy.

434 Figure 4a shows the resulting KGE values for the evaluated regionalization methods and the calibrated version as
435 grouped boxplots for different ranges of calibrated γ . The methods show different performances for different γ
436 ranges, indicating their strengths and weaknesses. For the smallest γ range, "0.1-0.2", the selected methods that
437 perform well during the split-sample test outperform the benchmark-to-beat. The better result for minimal γ ranges
438 is probably partially related to the advantage of the tuning, which leads to more predictions of 0.1 within the

439 regionalization. The benchmark-to-beat shows the best performance for γ values between 0.2 and 0.5. The good
 440 performance for basins with calibrated γ values between 0.2 and 0.5 is probably related to the benefit of using the
 441 logarithmical version of γ in the benchmark-to-beat, leading to more estimates of smaller values. However, this
 442 affects only 12 % of the basins, as calibrated values between 0.2 and 0.5 are not frequently present in the calibration
 443 result. Generally, the differences in KGE appear higher for smaller γ values, probably due to the decreasing pa-
 444 rameter sensitivity with higher values (see Appendix B).

445 Given the variability in the performance of the regionalization methods across the depicted γ ranges, it is challeng-
 446 ing to identify an overall best regionalization method using Fig. 4a. Therefore, we compare the various metrics of
 447 the KGE values depicted in Fig. 4b. The analyzed metrics are the minimum, maximum, mean, and median. Further,
 448 we count the number of poorly performing basins, defined as basins with a KGE below 0.2. In Fig. 4b, metrics
 449 that exceed the benchmark-to-beat are grey-shaded.



450 **Figure 4: a) KGE values of pseudo-ungauged basins from split-sample test grouped by the range of calibrated γ values,**
 451 **b) selected metrics of KGE values from the pseudo-ungauged basins (better or equal performance to the benchmark-**
 452 **to-beat is highlighted in grey), and c) histogram of the number of pseudo-ungauged basins with a KGE below 0.2 and**
 453 **the corresponding number of methods exhibiting this performance loss.**

454 Comparing the KGE metrics in Fig. 4b reveals that the methods showing higher logMAE values in our split-
 455 sampling test ensemble also show lower performance in simulating discharge. For example, all mean (and median)
 456 KGE values of the "worst" methods are below the mean KGE of 0.587 from the benchmark-to-beat, ranging from
 457 0.545 to 0.578. This indicates that the used logMAE between regionalized and calibrated values is a valid tool for
 458 a preliminary selection of adequate methods for the regionalization of WaterGAP3. However, for a more compre-
 459 hensive analysis, we recommend additionally analyzing the accuracy of simulated discharges, as the logMAE of
 460 calibrated and regionalized parameter values simplifies the inherent complexity between model parameters and
 461 model performance.

462 Moreover, SI (best) outperforms the benchmark-to-beat in all listed metrics, reducing poorly performing basins
 463 and enhancing well-performing basins. MLR (best) performs very similarly to SI (best), yet it shows a higher
 464 number of basins with KGE values below 0.2. In comparison to the benchmark-to-beat, it outperforms four out of

465 five criteria. The remaining well-performing methods, SP and knn (best), demonstrate superior or equal perfor-
466 mance to the benchmark-to-beat in three out of five criteria. SP results in an equal number of poorly performing
467 basins, and the minimal KGE value is lower than for the benchmark-to-beat. The knn (best) approach has a slightly
468 worse median of KGE, i.e., -0.001, and one additional basin shows a KGE below 0.2.

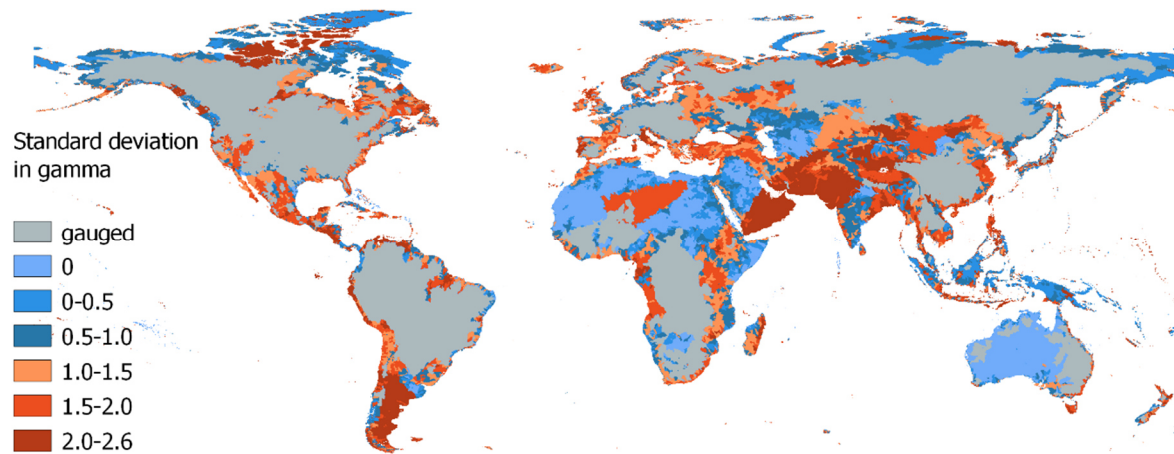
469 As SI (best) outperforms the benchmark-to-beat in all metrics, we conduct a statistical test to ascertain whether
470 there is a statistically significant difference in KGE results between the methods. To this end, we use a paired
471 Wilcoxon rank sum test to test the null hypothesis of whether the KGE differs significantly in central tendency. A
472 significance level of 0.05 and an adjusted p-value are applied to correct for multiple comparisons (using the cor-
473 rection after Benjamini & Hochberg (1995)). The results demonstrate that SI (best) outperforms all "worst" meth-
474 ods and the benchmark-to-beat. However, the null hypothesis for SP and the "best" options of knn and MLR cannot
475 be rejected. Consequently, rather than identifying a single alternative to the benchmark-to-beat, we have identified
476 four.

477 Notably, all regionalization methods lead to poorly performing basins, as evidenced by the range of basins with a
478 KGE below 0.2, varying from 13 to 37. In Fig. 4c, we examine whether there are basins that all methods cannot
479 regionalize, thereby indicating a general insufficiency of the regionalization methods for these basins. The histo-
480 gram indicates that most poorly performing basins belong to a single regionalization method. The high number of
481 basins, which cannot be estimated well by a single regionalization method, illustrates the diverse shortcomings of
482 the methods. A single basin shows poor performance across all methods. This is a basin of the river El Platanito
483 in Mexico. The calibrated γ value is about 1.5, and the corresponding KGE value in calibration is 0.466. This basin
484 appears to be highly sensitive to γ , with an inaccuracy in the estimated γ having a significant impact on the accuracy
485 of river discharge. For example, the benchmark-to-beat estimates γ to 1.0, which is close to the calibrated value of
486 1.5. However, the KGE value of the simulated discharge using the benchmark-to-beat is -0.158 due to a high
487 overestimation of the variation and mean of the discharge. This high sensitivity seems outstanding and is likely
488 attributable to the absence of waterbodies and snow, supporting a potentially high impact of γ on the model simu-
489 lation (Kupzig et al., 2023) in conjunction with a relatively small basin size (ca. 6,600 km²).

493 **3.4 Impacts on runoff simulations**

494 To evaluate the impact of runoff simulations, we apply an ensemble of regionalization methods generating γ esti-
495 mates for the worldwide ungauged regions. Within the ensemble, we use the four methods SI (best), knn (best),
496 MLR (best), and SP that (1) outperform the benchmark-to-beat regarding the logMAE of regionalized and cali-
497 brated values and (2) perform similarly to each other and better than the benchmark-to-beat in KGE for monthly
498 discharge. Additionally, we use the benchmark-to-beat as the fifth member of our regionalization method ensem-
499 ble. The entire set of 933 gauged basins is used for regionalizing γ , resulting in five distinct worldwide distributions
500 of γ . The spatially distributed standard deviation of the regionalized values is shown in Fig. 5.

501 In particular, the southern parts of South America, the northern and southern parts of North America, and Central
502 Asia reveal differences in γ across the ensemble of regionalization methods (see Fig. 5). In Europe, the highest
503 differences in regionalized values are observed in Italy, Great Britain, and northern Portugal. In Oceania, the high-
504 est values in standard deviation of γ are in Tasmania, New Zealand, and the southwest of Australia's coast. In
505 contrast, a minor variation in γ is apparent in northern Africa, most parts of Australia, and the East of the Dead
506 Sea. Thus, the uncertainty associated with globally regionalizing γ seem to vary across different regions.



507
 508 **Figure 5: Standard deviation in regionalized γ values using the best approaches of MLR (best), SI (best), SP, knn (best),**
 509 **and the benchmark-to-beat. Note that dry regions without discharge are set to zero.**

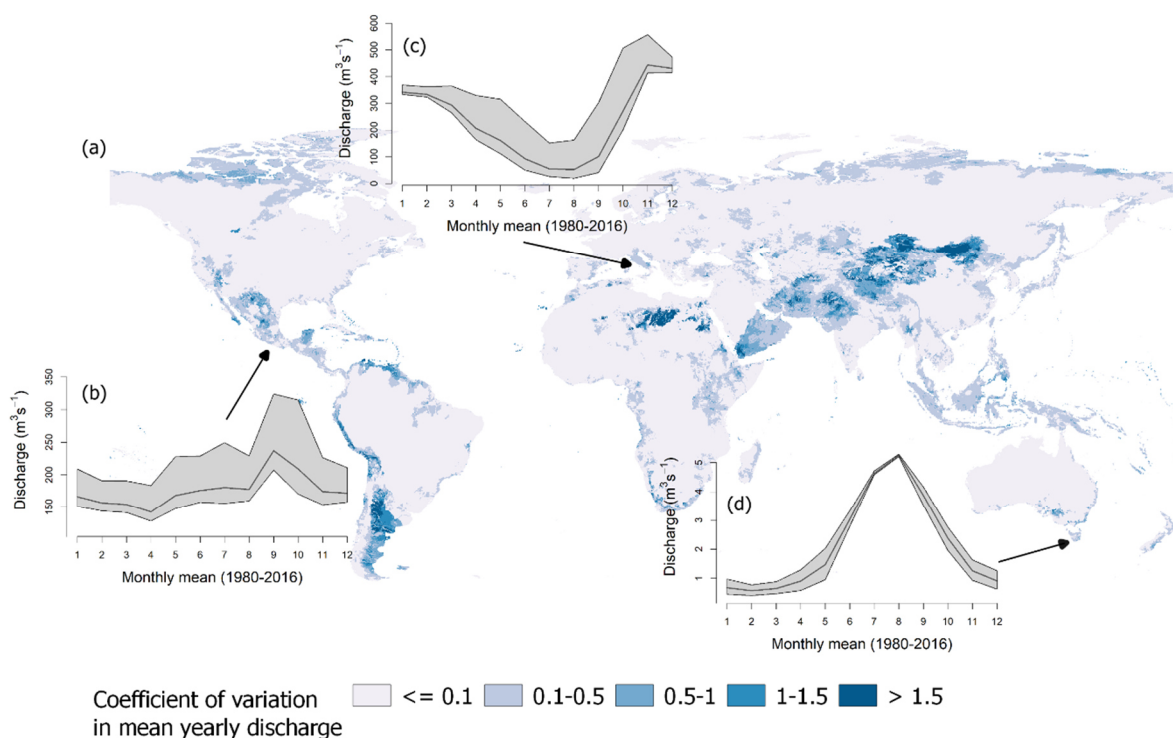
510 An example of how these uncertainties in regionalized values propagate through the water system is presented in
 511 Fig. 6. This figure displays the coefficient of variation of the mean yearly discharge between 1980 and 2016 based
 512 on the five simulation runs. Moreover, we highlight the effect on rivers in ungauged regions by showing the re-
 513 sulting seasonal pattern, i.e., the simulated long-term mean of monthly river discharge for three exemplary rivers.
 514 These rivers are the Río Bravo in Mexico, the Tiber in Italy, and the Tamar River in Tasmania. Each river is located
 515 in an ungauged region, where the standard deviation in γ is high (see Fig. 5).

516 Comparing Fig. 5 and Fig. 6 reveals that regions showing variability in γ tend to exhibit variation in mean yearly
 517 discharge. However, the impact of variation in γ on the simulated discharge appears to vary spatially. Some regions
 518 showing a high degree of variation in γ do not exhibit a correspondingly high degree of variation in discharge. For
 519 example, 45 % of all ungauged regions showing a low variation in discharge, i.e., the coefficient of variation is
 520 below 0.5, exhibit a standard deviation of more than one in γ . In contrast, about 89 % of the ungauged regions
 521 showing a higher discharge variation exhibit a standard deviation of more than one in γ . Thus, variation in γ does
 522 not necessarily lead to variation in river discharge, but it increases the likelihood that a region's discharge is af-
 523 fected. The spatially varying impact of γ is likely related to varying sensitivity regarding γ in the ungauged regions,
 524 which depends on numerous aspects, e.g., snow occurrence or waterbodies (see Kupzig et al., 2023).

525 About 11 % of the ungauged area exhibits variations in yearly river discharge exceeding 50 % of the mean. These
 526 regions are primarily in southern South America and Central Asia. A further 62 % of the ungauged area exhibits
 527 variations in yearly river discharge between 10 % and 50 % of the mean. These regions are mainly located on the
 528 northern coast of Russia and northern Canada, Indonesia, and Tasmania. Other areas, like most ungauged regions
 529 of Africa and Australia, show almost no impact, i.e., the variation in yearly discharge is less than 10 % of the
 530 mean. In northern Africa, one region exhibits higher values in the coefficients of variation. These values are at-
 531 tributable to minimal discharge values, resulting in comparatively high coefficients of variation in this region.

532 Considering the variation in the seasonality in the selected ungauged river systems (see Fig. 6b-d), the temporal
 533 impact of regionalization varies across the local landscape. For the Tamar River in Tasmania, as illustrated in Fig.
 534 6d, the variation is higher at the start and end of the dry periods in October/November and April/May, respectively.
 535 The spread in monthly mean discharge is about $0.7 \text{ m}^3\text{s}^{-1}$ to $1 \text{ m}^3\text{s}^{-1}$ in these periods. The Tiber in Italy and the Río
 536 Bravo in Mexico exhibit a similar pattern: using the regionalized γ values of SP leads to much higher discharge

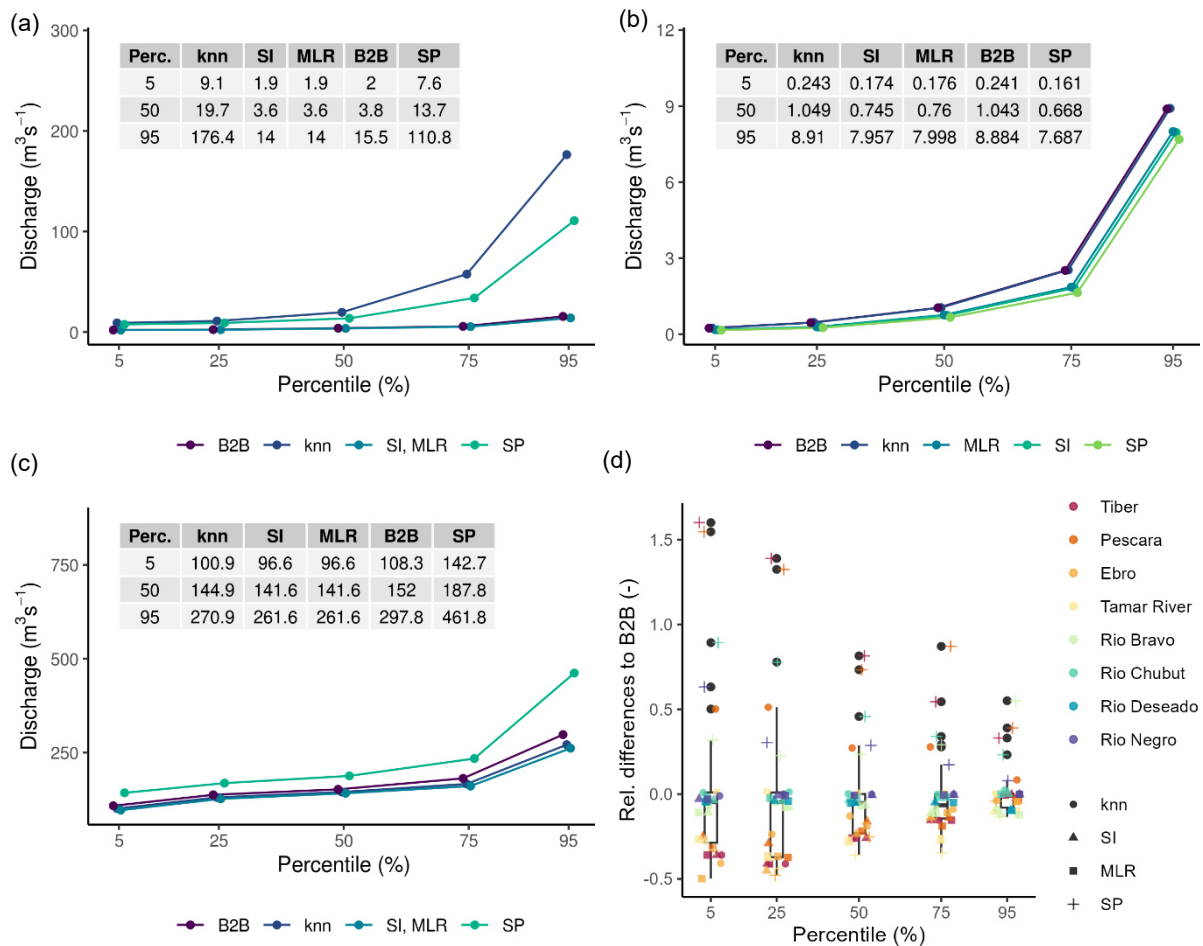
537 rates than other ensemble members, introducing broad uncertainty bands. For the Tiber, this leads to seasonal
 538 estimates varying between 1.2 % (in January) and 11 % (in October) of the mean yearly sum. The Río Bravo shows
 539 variations in its seasonal pattern, with values ranging from 2.2 % (in February) to 6.8 % (in October) of the mean
 540 yearly sum. Thus, all rivers display a temporally varying impact. Whereas the main variation in the discharge of
 541 the Río Bravo and the Tiber is mainly attributed to the SP regionalization run, for the Tamaris River, all regional-
 542 ization runs contribute to the varying long-term monthly mean in discharge.
 543



545 **Figure 6: a) Global map of the coefficient of variation in mean yearly discharge for the applied regionalization methods.**
 546 **Resulting differences in the regionalization ensemble regarding the long-term mean of monthly discharge are depicted**
 547 **for: b) the Río Bravo in Mexico, c) the Tiber in Italy and d) the Tamar River in Tasmania. The grey-shaded area**
 548 **indicates the range of the long-term mean of monthly discharge and the black line indicates the mean off all simulation**
 549 **runs.**

550 To gain a deeper understanding of the local impact of regionalization on runoff simulations, we analyze the annual
 551 percentiles from 1980 to 2016 for Río Deseado in Argentina, Río Bravo, and Tamar River, displaying the mean
 552 percentile of all years (see Fig. 7a-c). As the Tiber and Río Bravo display high similarities in the resulting patterns
 553 of percentiles, we demonstrate the impact by showing the percentiles from the Río Bravo. Additionally, we com-
 554 pare the relative differences in the mean for each percentile using eight ungauged river systems (see Fig. 7d), as
 555 previously done by Gudmundsson et al. (2012) for nine GHMs. To calculate the relative difference, we subtract
 556 the mean annual percentile of a method from the corresponding mean annual percentile of the reference and divide
 557 the resulting difference by the mean annual percentile of the reference. Instead of using observed flow as a refer-
 558 ence, we use the annual percentiles of our benchmark-to-beat. As river discharge is already spatially aggregated
 559 information, it is unnecessary to spatially aggregate grid cells to create results comparable to those of Gudmunds-
 560 son et al. (2012), who used cell runoff. The evaluated river systems are Río Chubut, Río Deseado, Río Negro, Río
 561 Bravo, Tamar River, Tiber, Pescara, and Ebro.

562



563 **Figure 7: Mean annual percentiles between 1980 and 2016 of simulated discharge using an ensemble of regionalization**
 564 **methods. The river are a) Río Deseado, b) Tamar River, and c) Río Bravo. In d), the relative differences in mean annual**
 565 **percentiles to the benchmark-to-beat of eight ungauged river systems are presented. Negative values indicate smaller**
 566 **mean annual percentiles than the benchmark-to-beat. Note that all data points from Río Deseado for knn and SP are**
 567 **excluded as the values are above 2.0.**

568 In Fig. 7a, Río Deseado is highly affected by uncertainties in simulated discharge due to the different regionaliza-
 569 tion methods; all segments of the percentiles show high variations where the absolute spread is increasing with
 570 increasing percentiles. For SP and knn (best), the discharge is highest, e.g., estimating a median discharge of 13.7
 571 m³s⁻¹ and 19.7 m³s⁻¹, respectively. For the other methods, the simulated discharge is low, e.g., SI and MLR result
 572 in an equal median discharge of 3.6 m³s⁻¹. The Tamar River in Fig. 7b also shows increasing absolute differences
 573 between the methods for higher percentiles, with the benchmark-to-beat approach leading to the highest discharge.
 574 For the Río Bravo, the absolute differences between the highest result of SP and the other methods remain almost
 575 constant until the 75th percentile. For the 95th percentile, the absolute differences increase rapidly from about 40
 576 m³s⁻¹ (75th percentile) to nearly 200 m³s⁻¹ (95th percentile). The exemplary results of Río Deseado and Río Bravo
 577 indicate a potentially high degree of uncertainty regarding the high percentiles in discharge simulation. These
 578 uncertainties put the results of global flood frequency analysis (e.g., Ward et al., 2013) in ungauged regions at risk
 579 as the time series of annual maxima might be even more uncertain. Thus, the results of flood frequency analysis
 580 should be carefully interpreted in ungauged regions as the impact of parameter regionalization may be significant.

581 Upon examination of the relative differences to the benchmark-to-beat for eight ungauged river systems, it be-
 582 comes evident that the impact of regionalization methods varies between ungauged river systems (e.g., Río Negro
 583 exhibits almost no variation, but Ebro does). Moreover, it becomes apparent that some regionalization methods

584 contribute more to the variation in estimated discharge than others. The methods contributing most are knn (best)
 585 and SP. For knn (best), 10 of the 40 relative differences are higher than |0.3|. For SP, even 29 out of the 40 relative
 586 differences are higher than |0.3|. The results of SI (best) and MLR (best) are very similar, indicating high similarity
 587 in performance. This is consistent with the KGE evaluation (see Chapter 3.3), in which they performed similarly.
 588 The observation in Fig. 7d that higher relative differences of discharge simulations occur in drier percentiles is
 589 also reported in Gudmundsson et al. (2012). Moreover, the relative differences between the five regionalization
 590 runs seem comparable to the inter-model differences depicted in Gudmundsson et al. (2012), indicating the high
 591 impact of regionalization methods on the evaluated ungauged river systems.

592 Finally, Table 3 presents the estimated yearly mean runoff to the ocean for all five ensemble members. All esti-
 593 mates of global "runoff to ocean" range from 45,622 (SI (best)) to 47,069 (SP). Thus, the differences are on the
 594 scale of smaller inter-model differences (see Table 2 in Widen-Nilsson et al.,2007). The impact of regionalization
 595 becomes even more evident using an unsuitable regionalization method for WaterGAP3. For instance, the tuned
 596 kmeans ("subset") approach results in 42,862 km³ yr⁻¹ "runoff to ocean", increasing the spread between the meth-
 597 ods to 4,208 km³ yr⁻¹ being in the scale of inter-model differences. This high impact of regionalization on global
 598 "runoff to ocean" is surprising, given that only 27 % of the world is ungauged, using the GRDC database. From
 599 this 27 %, most regions are in Australia and Africa, where minimal runoff is produced. In studies employing
 600 disparate models, e.g., for inter-model comparison, all regions are simulated in disparate ways.

601 The most significant deviations in the continental sums of "runoff to ocean" in Table 3 are due to SP. Only for
 602 Europe is the highest deviation related to MLR (best), not SP. Interestingly, the estimated sums of SP occasionally
 603 define the lowest and occasionally the highest extremes for the continents, lacking a systematic pattern. The out-
 604 standing role of SP is consistent with previous evaluations in this Chapter, where SP frequently contributes most
 605 to the variation in discharge. This suggests that SP may not be suitable for the global scale. Nevertheless, the
 606 pseudo-ungauged basins in the split-sample tests may also exhibit considerable distances from the observed basins.
 607 Given that SP achieved satisfactory results in both evaluations, using either the logMAE or the KGE, the evaluation
 608 indicates the method's suitability on a global scale. Thus, in the future, the split-sample test must be extended to
 609 gain deeper insights into the method's robustness and make a definitive statement about the method's suitability
 610 on a global scale. For example, the so-called "HDes" approach, recommended by Lebecherel et al. (2016), could
 611 be applied for this purpose. In this approach, the closest basin to the corresponding (pseudo-) ungauged basin is
 612 excluded from the regionalization process, thereby enabling an assessment of the method's robustness.

613 **Table 3: Mean outflow to the ocean and endorheic basins in km³ yr⁻¹ between 1980-2016. The highest continental devi-**
 614 **ation to the benchmark-to-beat is indicated in bold.**

<i>Runoff to ocean</i> ¹	B2B	SI (best)	knn (best)	MLR (best)	SP
Oceania	1,127	-1.80 %	-2.20 %	-3.40 %	-6.60 %
Europe	3,098	-2.30 %	-0.10 %	-2.60 %	0.20%
Asia	16,676	3.50 %	0.30 %	1.60 %	5.50 %
Africa	5,203	-1.00 %	0.70 %	-0.30 %	-3.60 %
North America	7,517	0.30 %	1.00 %	-1.70 %	2.20 %
South America	12,032	1.30 %	1.40 %	-0.20 %	4.90 %
global	45,653	46,273	45,953	45,622	47,069

¹including endorheic basin

615 **Conclusion**

616 Valid simulation results from GHMs, such as WaterGAP3, are crucial for detecting hotspots or studying patterns
617 in climate change impacts. However, the lack of worldwide monitoring data makes adapting GHMs' parameters
618 for valid global simulations challenging. Therefore, regionalization is necessary to estimate parameters in un-
619 gauged basins. This study applies regionalization methods for the first time to WaterGAP3, aiming to provide
620 insights into selecting suitable regionalization methods and evaluating their impact on the runoff simulations. Tra-
621 ditional and machine learning-based methods are tested to assess the application of several regionalization tech-
622 niques on a global scale. The concept of benchmark-to-beat and an ensemble of split-sampling tests are employed
623 for a comprehensive evaluation. Moreover, the impact on runoff simulation is assessed using a wide range of
624 temporal and spatial scales, i.e., from the daily to the yearly and from the local to the global scale.

625 In this study, four regionalization methods outperform the benchmark-to-beat and thus are considered appropriate
626 for WaterGAP3. These methods span the complete range of methodologies, i.e., regression-based methods and
627 methods using the concept of physical similarity and spatial proximity. Moreover, the methods vary in the de-
628 scriptors used to achieve optimal results. This highlights that different methods use descriptor sets with varying
629 efficiency. All methods perform best when using climatic and physiographic descriptors, indicating that combining
630 climatic and physiographic descriptors is optimal for regionalizing worldwide basins. Although random forest is
631 known to be especially robust among other machine learning-based techniques, it shows symptoms of over-pa-
632 rameterization, indicating that the algorithm is too flexible and adjusts to noise in the data, missing the underlying
633 systematic pattern.

634 Our results demonstrate that variation in the regionalized parameter value does not necessarily lead to variation in
635 river discharge. However, it increases the likelihood that a region's runoff is affected. This spatially varying impact
636 of γ is likely related to the varying sensitivity in ungauged regions regarding γ . Southern South America is a region
637 identified to be especially sensitive to variation in γ . Furthermore, local effects on runoff simulations indicate a
638 temporally varying impact. For example, some impacted rivers indicate a high degree of uncertainty regarding the
639 high percentiles in discharge simulation. These uncertainties potentially lead to a significant impact on flood fre-
640 quency analysis on a global scale, where the lack of gauging stations in certain regions calls for regionalization.
641 The global impact of regionalization methods that perform well for WaterGAP3 appears to be in the order of minor
642 inter-model differences. This impact rigorously increases when using a poorly performing method for WaterGAP3,
643 underscoring the importance of carefully selecting regionalization methods.

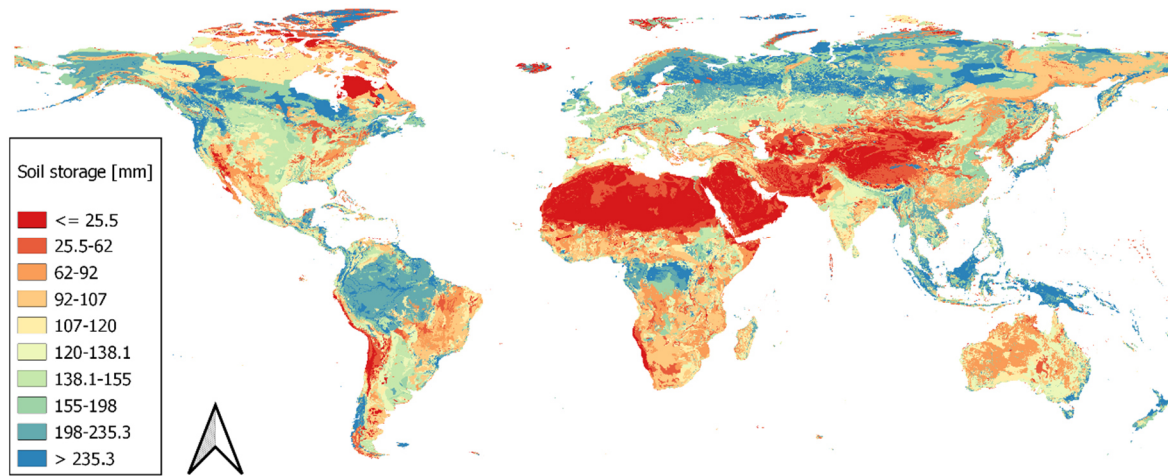
644 The spatial proximity approach contributes most to the variation in estimated runoff. The outstanding role of this
645 approach suggests that it may not be suitable for the global scale. However, as the pseudo-ungauged basins in the
646 split-sample tests may also have considerable large distances to the observed basins, and the method achieves
647 satisfactory results in all executed evaluations, it is not possible to make a definite statement about the method's
648 suitability for the global scale. Further research is required to gain deeper insights into the methods' robustness,
649 e.g., by extending the analysis by applying the recommended "HDes" approach (Lebecherel et al., 2016).

650 *Code and data availability.* The data and the supporting R-Code to reproduce this study's findings are available at
651 <https://doi.org/10.5281/zenodo.11833447>.

652 *Authors contribution.* JK developed, designed, and drafted the study. NK helped to design the experiment. MF
653 provided feedback throughout the entire process and supported the writing.

654 *Competing interests.* The authors declare that they have no conflict of interest.

655 **Appendix A: Global Map of derived global soil moisture storage**



656

657 **Figure A1: Global map of the size of soil storage based on Batjes (2012) and land use information (derived from Friedl**
658 **& Sulla-Menashe, 2019)**

659

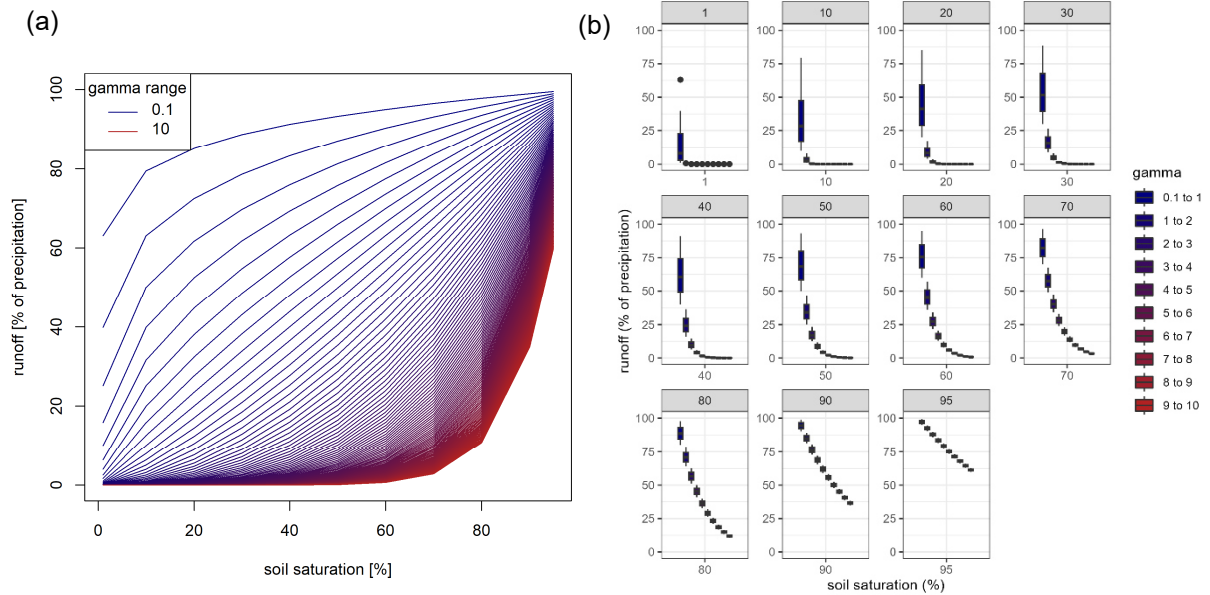
660 **Appendix B: Further analysis regarding the clustering of parameter values at the extremes**

661 The clustered calibrated parameter values at the extremes of the valid parameter space (see Fig. 1b) are a known
 662 problem within the calibration. As the parameter space, i.e., the parameter bounds, is crucial for calibration and,
 663 in consequence, for regionalization, we address this issue by a brief sensitivity analysis to demonstrate that the
 664 clustering of the calibrated parameter values is more an issue of missing processes (or using additional parameter
 665 values) than an issue of inappropriate parameter space. As the lower limit of the calibrated parameter (0.1) is
 666 sufficiently small in comparison to other studies using a similar HBV-based approach for runoff generation pro-
 667 cesses (e.g., see the beta in Table A2 in Jansen et al., 2022), we focus on the sensitivity analysis on the upper limit
 668 of γ (5.0).

669 In the sensitivity analysis regarding the upper limit of γ , we applied the model formula (see equation B1) containing
 670 the model's parameter γ and modified it within the bounds of 0.1 and 10. Additionally, we modified the soil satu-
 671 ration varying from 1 % to 95 %.

$$outflow = precipitation_{effective} \cdot soil\ saturation^{\gamma} \quad (B1)$$

672 The calculated outflow and its relationship to the soil saturation and γ are depicted in Fig. B1 and B2. The incoming
 673 effective precipitation is defined as constant. As it is a factor in equation B1,, the results regarding incoming
 674 effective precipitation are linearly scalable.



675 **Figure B1: a) Runoff generation in the soil layer (neglecting overflow and evapotranspiration) using different values**
 676 **for the calibration parameter and increasing the soil-moisture, b) runoff generation for varying soil moisture grouped**
 677 **in bins of size one.**

678 In the depicted Fig. B1, the runoff generation process differences between differing γ values become more linear
 679 when soil saturation increases. Thus, the non-linear model parameter becomes less critical for high soil moisture.
 680 Generally, the runoff generation process differences for higher γ values are more pronounced for higher soil mois-
 681 ture. For lower soil moisture, the smaller values have higher effects on the generated runoff. For example, for 70 %
 682 soil moisture, the differences for γ values ranging from 5 to 10 are between 3 % and 16 %. For the same soil
 683 moisture, the range in runoff generation varies from 16 % to 70 % for γ values between 1 and 5.

684 High γ values usually occur in dry regions (see Fig. 4b in Müller Schmied et al., 2021). In dry regions, high soil
685 moisture values are not expected to occur frequently (e.g., see Khosa et al., 2020; Oloruntoba et al., 2024 for
686 estimated and measured soil moisture in Africa and Draper et al., 2008 for estimated and measured soil moisture
687 in Australia). It is, therefore, unlikely that higher γ values will significantly enhance the calibration result or de-
688 crease the issue of clustered calibrated parameter values at the higher end of the parameter space. More likely, the
689 clustering of calibrated parameter values will be resolved in dry regions by incorporating additional (missing)
690 model processes, such as evaporation from rivers or inaccurate representation of groundwater processes (Eisner,
691 2016, p. 49). Thus, the parameter bounds of γ (e.g., also used in Eisner 2016, p. 16; Müller Schmied et al., 2021;
692 Müller Schmied et al., 2023) are not changed in this study.

693 **Appendix C: Basin descriptors**

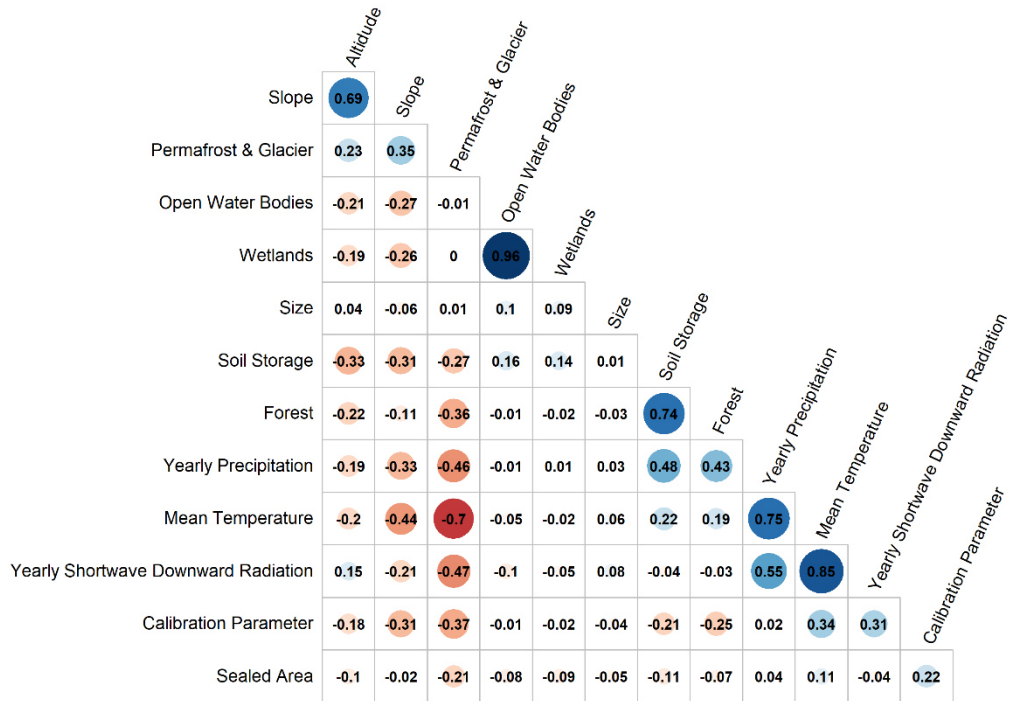
694 Overview of basins descriptors used in this study. All basin descriptors are derived from the original model input
695 and aggregated with a simple mean method to basin values to produce the same spatial resolution as the calibrated
696 model parameter.

- 697 • *Soil Storage*: The size of the soil storage, i.e., the maximal water content in the soil reachable for plants
698 in mm. The information is the product of rooting depth (defined in a look-up table) and the total available
699 water content derived from Batjes (2012).
- 700 • *Open Water Bodies*: The fraction of the area covered with open water bodies in the basin is given as a
701 percentage. The model input is based on the GLWD database (Lehner & Döll, 2004).
- 702 • *Wetlands*: The fraction of area covered with wetlands in a basin is given in percentage. The model input
703 is based on the GLWD database (Lehner & Döll, 2004).
- 704 • *Size*: Size of a basin in km².
- 705 • *Slope*: The mean slope class is calculated as described in Döll & Fiedler (2008) and based on GTOPO30
706 (USGS EROS data centre).
- 707 • *Altitude*: The mean altitude of a basin is given in meters above sea level and based on GTOPO30 (USGS
708 EROS data centre).
- 709 • *Forest*: The mean fraction of the area covered with forest is given in percentage and derived from MODIS
710 data (Friedl & Sulla-Menashe, 2019), where 2001 is used as a reference. All grid cells having a dominant
711 International Geosphere-Biosphere Programme (IGBP) classification between one and five are defined
712 as "forest".
- 713 • *Sealed Area*: The mean fraction of sealed area is given in percentage and derived from MODIS data
714 (Friedl & Sulla-Menashe, 2019), where 2001 is used as a reference. All grid cells having an IGBP clas-
715 sification equal to 13 are defined as they would contain 60% of the sealed area. Note: The different treat-
716 ment of forest and sealed area is based on the required model input; whereas the land cover is a classified
717 value, the sealed area is a floating-point value.
- 718 • *Permafrost & Glacier*: The mean coverage of permafrost and glacier in a basin is given in percentage. It
719 is based on the World Glacier Inventory and the Circum-Arctic Map of Permafrost and Ground-Ice Con-
720 ditions.
- 721 • *Mean Temperature*: The mean air temperature is based on the meteorological forcing used to drive the
722 model (Lange, 2019) covering the period 1979 to 2016 and given in degrees Celsius.
- 723 • *Yearly Precipitation*: The yearly precipitation sum is based on the meteorological forcing used to drive
724 the model (Lange, 2019) covering the period 1979 to 2016 and given in mm.
- 725 • *Yearly Shortwave Downward Radiation*: The yearly shortwave downward radiation is based on the me-
726 teorological forcing used to drive the model (Lange, 2019) covering the period 1979 to 2016 and given
727 in Wm⁻².

728

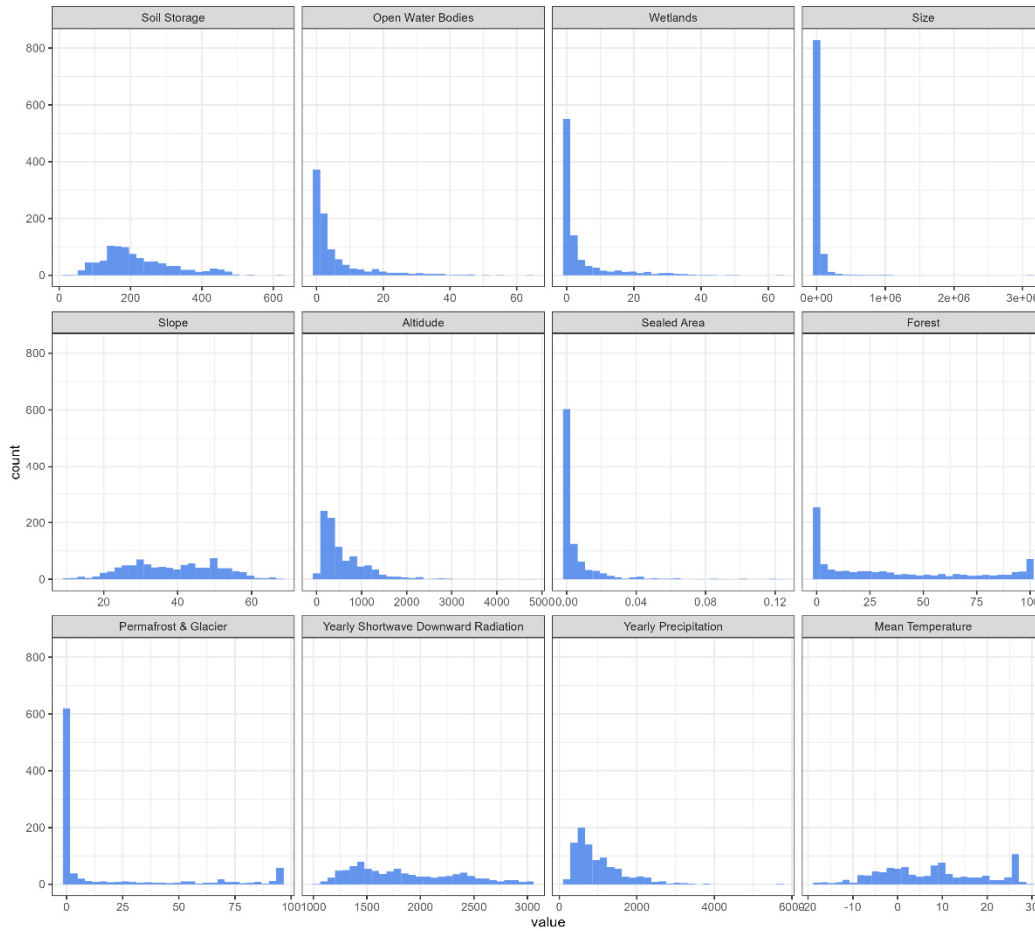
729 The correlation between the defined basin descriptors is shown in Fig. A1. The variation within each basin de-
730 scriptor for basins used for regionalization is shown in Fig. A2.

731



732
733
734
735

Figure C1: Correlation between basins descriptors.



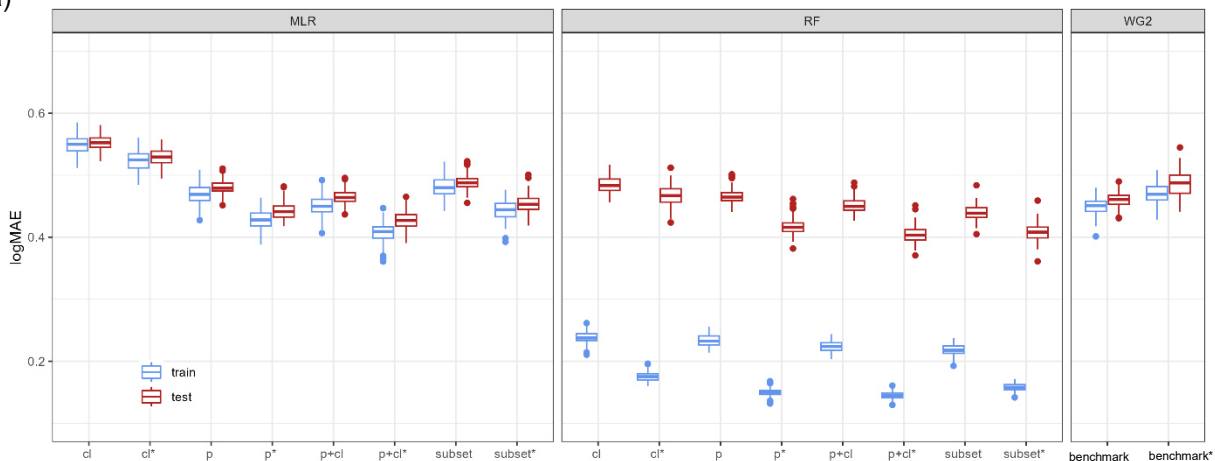
736
737

Figure C2: Distribution of basins descriptors within all basins used for regionalization (n=933)

738

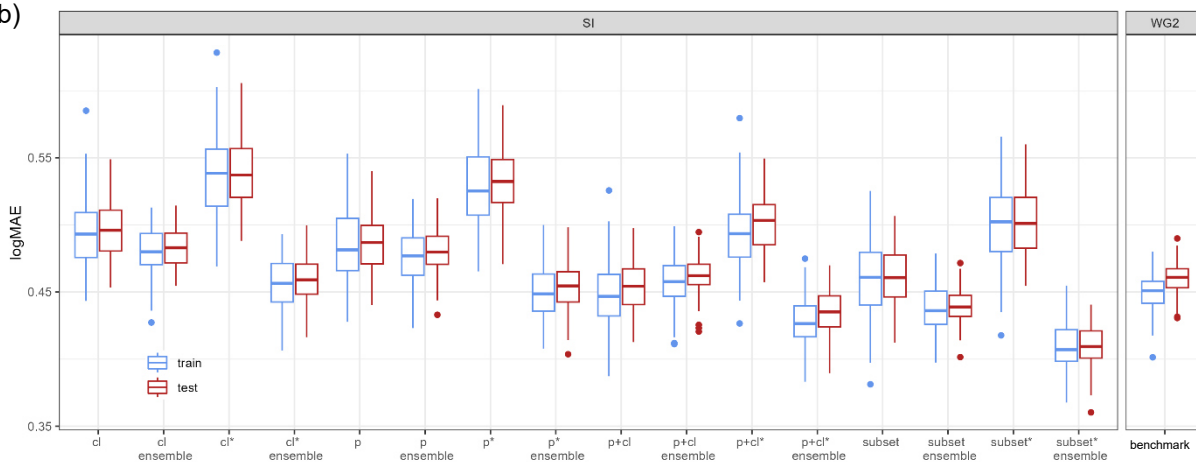
739 **Appendix D: Results of the ensemble of the split-sample tests**

(a)



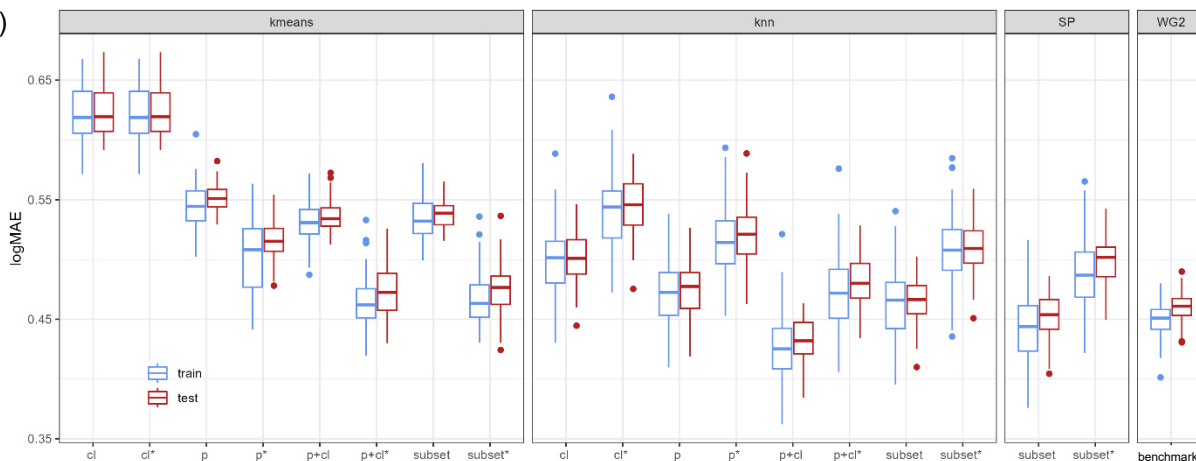
740

(b)



741

(c)



742

743 **Figure D1: logMAE values for all 100 split-sampling tests using all variants of a) MLR, RF, and benchmark-to-beat,**
744 **b) SI, and c) kmeans, knn, and SP. Note that the asterisk * indicates the tuned version of the method.**

745

746 **Table D1: Performance loss in median logMAE of the ensemble of split-sample tests from training to testing expressed**
 747 **in % of logMAE in training.**

test (% train)	MLR	RF	SI		kmeans	knn	SP	B2B
			no ens.	ensem- ble				
cl	100.4	202.9	100.6	100.6	100	100	102.3	102.2
p	102.1	199.6	101.2	100.6	101.3	101.1		
p+cl	103.1	207.1	101.6	100.9	100.6	95.6		
subset	101.7	223.9	100	100.7	101.3	100.2		

test* (% train*)	MLR	RF	SI		kmeans	knn	SP	B2B
			no ens.	ensem- ble				
cl	100.8	266.9	99.8	100.7	100	100.4	103.1	104.1
p	103	277.3	101.3	101.3	101.4	101.4		
p+cl	104.4	277.9	102	102.1	102.2	101.7		
subset	102	258.2	99.8	100.5	103	100.2		

748

749

750 **References**

- 751 Arheimer, B., Pimentel, R., Isberg, K., Crochemore, L., Andersson, J. C. M., Hasan, A., & Pineda, L.: Global
752 catchment modelling using World-Wide HYPE (WWH), open data, and stepwise parameter estimation, *Hydrology
753 and Earth System Sciences*, 24(2), 535–559. <https://doi.org/10.5194/hess-24-535-2020>, 2020.
- 754 Arsenault, R., & Brissette, F. P.: Continuous streamflow prediction in ungauged basins: The effects of equifinality
755 and parameter set selection on uncertainty in regionalization approaches, *Water Resources Research*, 50, 6135–
756 6153, <https://doi.org/10.1002/2013WR014898>, 2014.
- 757 Ayzel, G. V., Gusev, E. M., & Nasonova, O. N.: River runoff evaluation for ungauged watersheds by SWAP
758 model. 2. Application of methods of physiographic similarity and spatial geostatistics, *Water Resources*, 44(4),
759 547–558, <https://doi.org/10.1134/S0097807817040029>, 2017.
- 760 Barbarossa, V., Bosmans, J., Wanders, N., King, H., Bierkens, M. F. P., Huijbregts, M. A. J., & Schipper, A. M.:
761 Threats of global warming to the world's freshwater fishes, *Nature Communications*, 12(1), 1701,
762 <https://doi.org/10.1038/s41467-021-21655-w>, 2021.
- 763 Batjes, N. H.: ISRIC-WISE derived soil properties on a 5 by 5 arc-minutes global grid (ver. 1.2) [data set],
764 <https://data.isric.org/geonetwerk/srv/eng/catalog.search#/metadata/82f3d6b0-a045-4fe2-b960-6d05bc1f37c0>,
765 2012.
- 766 Beck, H. E., Pan, M., Lin, P., Seibert, J., van Dijk, A. I. J. M., & Wood, E. F: Global Fully Distributed Parameter
767 Regionalization Based on Observed Streamflow From 4,229 Headwater Catchments, *Journal of Geophysical Re-
768 search: Atmospheres*, 125(17), <https://doi.org/10.1029/2019JD031485>, 2020.
- 769 Beck, H. E., van Dijk, A. I. J. M., Roo, A. de, Dutra, E., Fink, G., Orth, R. & Schellekens, J.: Global evaluation of
770 runoff from 10 state-of-the-art hydrological models, *Hydrol. Earth Syst. Sci.*, 21, 2881–20903,
771 <https://doi.org/10.5194/hess-21-2881-2017>, 2017.
- 772 Beck, H. E., van Dijk, A. I. J. M., Roo, A. de, Miralles, D. G., McVicar, T. R., Schellekens, J., & Bruijnzeel, L.
773 A.: Global-scale regionalization of hydrologic model parameters, *Water Resources Research*, 52(5), 3599–3622,
774 <https://doi.org/10.1002/2015WR018247>, 2016.
- 775 Benjamini, Y., & Hochberg, Y: Controlling the False Discovery Rate: A Practical and Powerful Approach to
776 Multiple Testing, *Journal of the Royal Statistical Society. Series B (Methodological)*, 57(1), 289–300.
777 <http://www.jstor.org/stable/2346101>, 1995.
- 778 Boulange, J, Hanasaki, N, Yamazaki, D., & Pokhrel, Y.: Role of dams in reducing global flood exposure under
779 climate change, *Nature Communications*, 12(1), 417, <https://doi.org/10.1038/s41467-020-20704-0>, 2021.
- 780 Breimann, L.: Random Forests, *Machine Learning*, 45, 1–32, <https://doi.org/10.1023/A:1010933404324>, 2001.
- 781 Chaney, N. W., Herman, J. D., Ek, M. B., & Wood, E. F.: Deriving global parameter estimates for the Noah land
782 surface model using FLUXNET and machine learning, *Journal of Geophysical Research: Atmospheres*, 121(22),
783 13,218–13,235, <https://doi.org/10.1002/2016JD024821>, 2016.
- 784 Charrad, M., Ghazzali, N., Boiteau, V., Niknafs, A.: NbClust: An R Package for Determining the Relevant Number
785 of Clusters in a Data Set, *Journal of Statistical Software*, 61(6), 1–36. <https://doi.org/10.18637/jss.v061.i06>, 2014.

786 Cuntz, M., Mai, J., Samaniego, L., Clark, M., Wulfmeyer, V., Branch, O., Attinger, S., & Thober, S.: The impact
787 of standard and hard-coded parameters on the hydrologic fluxes in the Noah-MP land surface model, *Journal of*
788 *Geophysical Research: Atmospheres*, 121, 10,676 - 10,700, <https://doi.org/10.1002/2016JD025097>, 2016.

789 Döll, P. & Fiedler, K.: Global-scale modeling of groundwater recharge, *Hydrol. Earth Syst. Sci.*, 12, 863–885,
790 <https://doi.org/10.5194/hess-12-863-2008>, 2008

791 Döll, P., Kaspar, F., & Lehner, B.: A global hydrological model for deriving water availability indicators: model
792 tuning and validation, *Journal of Hydrology*, 270, 105–13, [https://doi.org/10.1016/S0022-1694\(02\)00283-4](https://doi.org/10.1016/S0022-1694(02)00283-4), 2003.

793 Döll, P., Hasan, H. M. M., Schulze, K., Gerdener, H., Börger, L., Shadkam, S., Ackermann, S., Hosseini-Moghari,
794 S.-M., Müller Schmied, H., Güntner, A., & Kusche, J.: averaging multi-variable observations to reduce and quan-
795 tify the output uncertainty of a global hydrological model: evaluation of three ensemble-based approaches for the
796 Mississippi River basin, *Hydrology and Earth System Sciences*, 28 (10), 2259-2295, [https://doi.org/10.5194/hess-](https://doi.org/10.5194/hess-28-2259-2024)
797 [28-2259-2024](https://doi.org/10.5194/hess-28-2259-2024), 2024.

798 Draper, C. S., Walker, J. P., Steinle, P. J., de Jeu, R. A. M., Holmes T. R. H.: An evaluation of AMSR–E derived
799 soil moisture over Australia, *Remote Sensing of Environment*, 113, 703-710,
800 <https://doi.org/10.1016/j.rse.2008.11.011>, 2008.

801 Eisner, S.: Comprehensive Evaluation of the WaterGAP3 Model across Climatic, Physiographic, and Anthro-
802 pogenic Gradients, Ph.D. thesis, University of Kassel, Kassel, Germany, 128pp., 2016.

803 Friedl, M., Sulla-Menashe, D.: MCD12Q1 MODIS/Terra+Aqua Land, Cover Type Yearly L3 Global 500m SIN
804 Grid V006, NASA EOSDIS Land Processes DAAC [data set], <https://doi.org/10.5067/MODIS/MCD12Q1.006>,
805 2019.

806 Feigl, M., Thober, S., Schweppe, R., Herrnegger, M., Samaniego, L., & Schulz, K.: Automatic Regionalization of
807 Model Parameters for Hydrological Models, *Water Resources Research*, 58, e2022WR031966,
808 <https://doi.org/10.1029/2022WR031966>, 2022.

809 Golian, S., Murphy, C., & Meresa, H.: Regionalization of hydrological models for flow estimation in ungauged
810 catchments in Ireland, *Journal of Hydrology: Regional Studies*, 36, 100859,
811 <https://doi.org/10.1016/j.ejrh.2021.100859>, 2021.

812 GRDC, The Global Runoff Data Centre, 56068 Koblenz, Germany, 2020.

813 Gudmundsson, L., Tallaksen, L. M., Stahl, K., Clark, D. B., Dumont, E., Hagemann, S., Bertrand, N., Gerten, D.,
814 Heinke, J., Hanasaki, N., Voss, F., & Koirala, S.: Comparing Large-Scale Hydrological Model Simulations to
815 Observed Runoff Percentiles in Europe. *Journal of Hydrometeorology*, 13(2), 604-620.
816 <https://doi.org/10.1175/JHM-D-11-083.1>, 2012.

817 Guo Y, Zhang Y, Zhang L, & Wang Z: Regionalization of hydrological modeling for predicting streamflow in
818 ungauged catchments: A comprehensive review, *WIREs Water*, 8, e1487, <https://doi.org/10.1002/wat2.1487>,
819 2020.

820 Gupta, H. V, Sorooshian, S., & Yapo, P. O.: Toward improved calibration of hydrologic models: Multiple and
821 noncommensurable measures of information, *Water Resources Research*, 34(4), 751–763,
822 <https://doi.org/10.1029/97WR03495>, 1998.

823 He, Y., Bárdossy, A., & Zehe, E.: A review of regionalisation for continuous streamflow simulation, *Hydrology*
824 *and Earth System Sciences*, 15(11), 3539–3553. <https://doi.org/10.5194/hess-15-3539-2011>, 2011.

825 Jansen, K. F., Teuling, A. J., Craig, J. R., Dal Molin, M., Knoben, W. J. M., Parajka, J., Vis, M., Melsen, L. A.:
826 Mimicry of a conceptual hydrological model (HBV): What's in a name? *Water Resources Research*, 57,
827 e2020WR029143. <https://doi.org/10.1029/2020WR029143>, 2022.

828 Kaspar, F.: Entwicklung und Unsicherheitsanalyse eines globalen hydrologischen Modells, Ph.D. thesis, Univer-
829 sity of Kassel, Kassel, Germany, 129pp., 2004.

830 Khosa, F. V., Mateyisi, M. J., van der Merwe, M. R., Feig, G. T., Engelbrecht, F. A., Savage, M. J.: Evaluation of
831 soil moisture from CCAM-CABLE simulation, satellite-based models estimates and satellite observations: a case
832 study of Skukuza and Malopeni flux towers, *Hydrology and Earth System Sciences*, 24(4), 1587-1609,
833 <https://doi.org/10.5194/hess-24-1587-2020>, 2020.

834 Krabbenhoft, C. A., Allen, G. H., Lin, P., Godsey, S. E., Allen, D. C., Burrows, R. M., DelVecchia, A. G., Fritz,
835 K. M., Shanafield, M., Burgin, A. J., Zimmer, M. A., Datry, T., Dodds, W. K., Jones, C. N., Mims, M. C., Franklin,
836 C., Hammond, J. C., Zipper, S., Ward, A. S., Olden, J. D.: Assessing placement bias of the global river gauge
837 network, *Nature Sustainability*, 5, 586–592. <https://doi.org/10.1038/s41893-022-00873-0>, 2022.

838 Kupzig, J., Reinecke, R., Pianosi, F., Flörke, M., & Wagener, T.: Towards parameter estimation in global hydro-
839 logical models, *Environmental Research Letters*, 18(7), 74023. <https://doi.org/10.1088/1748-9326/acdae8>, 2023.

840 Lange, S.: Earth2Observe, WFDEI and ERA-Interim data Merged and Bias-corrected for ISIMIP (EWEMBI), V.
841 1.1, GFZ Data Services [data set] , <https://doi.org/10.5880/pik.2019.004>, 2019.

842 Lebecherel, L., Andréassian, V., Perrin: On evaluating the robustness of spatial-proximity-based regionalization
843 methods, *Journal of Hydrology*, 539, 196-203, <https://doi.org/10.1016/j.jhydrol.2016.05.031>, 2016.

844 Lehner, B. and Döll, P.: Development and validation of a global database of lakes, reservoirs and wetlands, *Journal*
845 *of Hydrology*, 296 (1-4), 1-22, <https://doi.org/10.1016/j.jhydrol.2004.03.028>, 2004.

846 Lehner, B., Verdin, K., & Jarvis, A.: New global hydrography derived from spaceborne elevation data, *Eos, Trans-*
847 *actions, AGU*, 89, 93–94, doi:10.1029/2008EO100001, 2008.

848 Liam, A., & Wiener, M.: Classification and Regression by randomForest. *R News*, 2(3), 18–22, 2002.

849 Lindström, G., Johansson, B., Persson, M., Gardelin, M., & Bergström, S.: Development and test of the distributed
850 HBV-96 hydrological model, *Journal of Hydrology*, 201, 272–288, <https://doi.org/10.1016/S0022->
851 [1694\(97\)00041-3](https://doi.org/10.1016/S0022-1694(97)00041-3), 1997.

852 McIntyre, N, Lee, H., Wheeler, H., Young, A., & Wagener, T.: Ensemble predictions of runoff in ungauged catch-
853 ments, *Water Resources Research*, 41(12), W12434, <https://doi.org/10.1029/2005WR004289>, 2005.

854 Merz, R., Blöschl, G.: Regionalisation of catchment model parameters, *Journal of Hydrology*, 287, 95-123,
855 <https://doi.org/10.1016/j.jhydrol.2003.09.028>, 2004.

856 Müller Schmied, H., Cáceres, D., Eisner, S., Flörke, M., Herbert, C., Niemann, C., Peiris, T. A., Popat, E., Port-
857 mann, F. T., Reinecke, R., Schumacher, M., Shadkam, S., Telteu, C.-E., Trautmann, T., Döll, P.: The global water

858 resources and use model WaterGAP v2.2d: model description and evaluation, *Geoscientific Model Development*,
859 14(2), 1037–1079, <https://doi.org/10.5194/gmd-14-1037-2021>, 2021.

860 Müller Schmied, H., Trautmann, T., Ackermann, S., Cáceres, D., Flörke, M., Gerdener, H., Kynast, E., Peiris, T.
861 A., Schiebener, L., Schumacher, M., Döll, P.: The global water resources and use model WaterGAP v2.2e: de-
862 scription and evaluation of modifications and new features, *Geoscientific Model Development Discussions* [pre-
863 print], 1-46, <https://doi.org/10.5194/gmd-2023-213>, 2023.

864 Nijssen, B., O'Donnell, G. M., Lettenmeier, D. P., Lohmann, D., & Wood, E. F.: Predicting the Discharge of
865 Global Rivers, *American Meteorological Society*, 3307–3323, [https://doi.org/10.1175/1520-0442\(2001\)014<3307:PTDOGR>2.0.CO;2](https://doi.org/10.1175/1520-0442(2001)014<3307:PTDOGR>2.0.CO;2), 2000.

867 Oloruntoba, B., Kollet, S., Motzka, C., Vereecken H., Franssen H.-J. H.: High Resolution Land Surface Modelling
868 over Africa: the role of uncertain soil properties in combination with temporal model resolution, *EGUsphere Pre-*
869 *print repository* [preprint], <https://doi.org/10.5194/egusphere-2023-3132>, 2024.

870 Oudin, L., Andréassian, V., Perrin, C., Michel, C., & Le Moine, N.: Spatial proximity, physical similarity, regres-
871 sion and ungauged catchments: A comparison of regionalization approaches based on 913 French catchments, *Water Resources Research*, 44(3), W03413, <https://doi.org/10.1029/2007WR006240>, 2008.

873 Oudin, L., Kay, A., Andréassian, V., & Perrin, C.: Are seemingly physically similar catchments truly hydrologi-
874 cally similar? *Water Resources Research*, 46(11), W11558, <https://doi.org/10.1029/2009WR008887>, 2010.

875 Pagliero, L., Bouraoui, F., Diels, J., Willems, P., & McIntyre, N.: Investigating regionalization techniques for
876 large-scale hydrological modelling, *Journal of Hydrology*, 570, 220–235, <https://doi.org/10.1016/j.jhydrol.2018.12.071>, 2019.

878 Parajka, J., Merz, R., & Blöschl, G.: A comparison of regionalisation methods for catchment model parameters,
879 *Hydrology and Earth System Sciences*, 9, 157–171, <https://doi.org/10.5194/hess-9-157-2005>, 2005.

880 Poissant, D., Arsenaault, R. & Brissette, F.: Impact of parameter set dimensionality and calibration procedures on
881 streamflow prediction at ungauged catchments, *Journal of Hydrology: Regional Studies*, 12, 220–237,
882 <https://doi.org/10.1016/j.ejrh.2017.05.005>, 2017.

883 Pool, S., Vis, M., & Seibert, J.: Regionalization for ungauged catchments — Lessons learned from a comparative
884 large-sample study. *Water Resources Research*, 57, e2021WR030437. <https://doi.org/10.1029/2021WR030437>,
885 2021.

886 Qi, W., Chen, J., Li, L., Xu, C., Li, J., Xiang, Y., & Zhang, S.: A framework to regionalize conceptual model
887 parameters for global hydrological modelling, *Hydrology and Earth System Sciences Discussions* [preprint],
888 <https://doi.org/10.5194/hess-2020-127>, 2020.

889 R Core Team.: R: A language and environment for statistical computing R Foundation for Statistical Computing,
890 Vienna, Austria. <https://www.r-project.org/>, 2020.

891 Reichl, J. P. C., Western, A. W., McIntyre, N. R. & Chiew, F. H. S: Optimization of a Similarity Measure for
892 Estimating Ungauged Streamflow, *Water Resources Research*, 45 (10), <https://doi.org/10.1029/2008WR007248>,
893 2009

894 Samaniego, L, Kumar, R & Attinger, S.: Multiscale parameter regionalization of a grid-based hydrologic model
895 at the mesoscale, *Water Resources Research*, 46(5), W05523, <https://doi.org/10.1029/2008WR007327>, 2010.

896 Schaeffli, B., & Gupta, H. V.: Do Nash values have value?, *Hydrological Processes*, 21(15), 2075–2080,
897 <https://doi.org/10.1002/hyp.6825>, 2007.

898 Schweppe, R., Thober, S., Müller, S., Kelbling, M., Kumar, R., Attinger, S., & Samaniego, L.: MPR 1.0: a stand-
899 alone multiscale parameter regionalization tool for improved parameter estimation of land surface models, *Geo-
900 scientific Model Development*, 15, 859–882, <https://doi.org/10.5194/gmd-15-859-2022>, 2022.

901 Seibert, J.: On the need for benchmarks in hydrological modelling, *Hydrological Processes*, 15(6), 1063–1064,
902 <https://doi.org/10.1002/hyp.446>, 2001.

903 Shannon, C. E.: A Mathematical Theory of Communication, *The Bell System Technical Journal*, 3(27), 379-423,
904 <https://doi.org/10.1002/j.1538-7305.1948.tb01338.x>, 1948.

905 Stacke, T., & Hagemann, S.: HydroPy (v1.0): a new global hydrological model written in Python, *Geoscientific
906 Model Development*, 14, 7795–7816, <https://doi.org/10.5194/gmd-14-7795-2021>, 2021.

907 Tang, Y., Marshall, L., Sharma, A. & Smith, T.: Tools for investigating the prior distribution in Bayesian hydro-
908 logy, *Journal of Hydrology*, 538, 551-562, <https://doi.org/10.1016/j.jhydrol.2016.04.032>, 2016.

909 Tongal, H., & Sivakumar, B.: Cross-entropy clustering framework for catchment classification, *Journal of Hydrol-
910 ogy*, 552, 433–446, <https://doi.org/10.1016/j.jhydrol.2017.07.005>, 2017.

911 Venables, W. N., & Ripley, B. D.: *Modern Applied Statistics with S (Fourth Edition)*. Springer Science+Business
912 Media New York, USA, 501pp, ISBN 978-1-4419-3008-8, 2002

913 Wagener, T., Wheeler, H. S., & Gupta, H. V.: *Rainfall – Runoff Modelling in Gauged and Ungauged Catchments*,
914 Imperial College Press, London, UK, 332pp., <https://doi.org/10.1142/p335>, 2004.

915 Wagener, T., & Wheeler, H. S.: Parameter estimation and regionalization for continuous rainfall-runoff models
916 including uncertainty, *Journal of Hydrology*, 320, 132-154, <https://doi.org/10.1016/j.jhydrol.2005.07.015>, 2006.

917 Ward, P. J., Jongman, B., Sperna Weiland, F., Bouwman, A., Van Beek, R., Bierkens, M. F. P., Ligtvoet, W., &
918 Winsemius, H. C.: Assessing flood risk at the global scale: model setup, results, and sensitivity, *Environmental
919 Research Letters*, 8, Article 044019. <https://doi.org/10.1088/1748-9326/8/4/044019>, 2013

920 Widén-Nilsson, E., Halldin, S., & Xu, C.: Global water-balance modelling with WASMOD-M: Parameter estimation and regionalisa-
921 tion, *Journal of Hydrology*, 340(1-2), 105–118, <https://doi.org/10.1016/j.jhydrol.2007.04.002>, 2007.

922 Wu, H., Zhang, J., Bao, Z., Wang, G., Wang, W., Yang, Y. & Wang, J.: Runoff Modeling in Ungauged Catchments
923 Using Machine Learning Algorithm-Based Model Parameters Regionalization Methodology, *Engineering*, 28, 93-
924 104, <https://doi.org/10.1016/j.eng.2021.12.014>, 2023.

925 Yang, X., Magnusson, J., Huang, S., Beldring, S., & Xu, C.: Dependence of regionalization methods on the com-
926 plexity of hydrological models in multiple climatic regions, *Journal of Hydrology*, 582, 124357,
927 <https://doi.org/10.1016/j.jhydrol.2019.124357>, 2020.

928 Yoshida, T., Hanasaki, N, Nishina, K., Boulange, J, Okada, M., & Troch, P. A.: Inference of Parameters for a
929 Global Hydrological Model: Identifiability and Predictive Uncertainties of Climate-Based Parameters, Water Re-
930 sources Research, 58, e2021WR030666, <https://doi.org/10.1029/2021WR030660>, 2022.