# Benchmarking multi-component spatial metrics for hydrologic model calibration using MODIS AET and LAI products

Eymen Berkay Yorulmaz<sup>1</sup>, Elif Kartal<sup>1</sup>, and Mehmet Cüneyd Demirel<sup>1</sup>

<sup>1</sup>Istanbul Technical University

August 24, 2023

#### Abstract

SPAtial EFficiency (SPAEF) metric is one of the most thoroughly metrics in hydrologic community. In this study, our aim is to improve SPAEF by replacing the histogram match component with other statistical indices, i.e. kurtosis and earth mover's distance, or by adding a fourth or fifth component such as kurtosis and skewness. The existing spatial metrics i.e. SPAtial Efficiency (SPAEF), Structural Similarity (SSIM) and Spatial Pattern Efficiency Metric (SPEM) were compared with newly proposed metrics to assess their converging performance. The mesoscale Hydrologic Model (mHM) of the Moselle River is used to simulate streamflow (Q) and actual evapotranspiration (AET). The two-source energy balance (TSEB) AET during the growing season is used as monthly reference maps to calculate the spatial performance of the model. The Moderate Resolution Imaging Spectroradiometer (MODIS) based Leaf area index (LAI) is utilized by the mHM via pedo-transfer functions and multi-scale parameter regionalization approach to scale the potential ET. In addition to the real monthly AET maps, we also tested these metrics using a synthetic true AET map simulated with a known parameter set for a randomly selected day. The results demonstrate that the newly developed four-component metric i.e. SPAtial Hybrid 4 (SPAH4) slightly outperform conventional three-component metric i.e. SPAEF (3% better). However, SPAH4 significantly outperforms the other existing metrics i.e. 40% better than SSIM and 50% better than SPEM. We believe that other fields such as remote sensing, change detection, function space optimization and image processing can also benefit from SPAH4.

# Benchmarking multi-component spatial metrics for hydrologic model calibration using MODIS AET and LAI products

# 3 E. B. Yorulmaz<sup>1</sup>, E. Kartal<sup>1</sup> and M. C. Demirel<sup>1</sup>

<sup>1</sup> Department of Civil Engineering, Istanbul Technical University, 34467 Maslak, Istanbul,
 Turkey

6

Corresponding author: Eymen Berkay Yorulmaz (<u>yorulmaz21@itu.edu.tr</u>, ORCID 0000-0003-3370-9465)

# 9 Key Points:

- Newly proposed spatial metrics offer significant improvements in discriminating
   between two raster maps
- Selecting appropriate spatial metric proved to be very crucial even for the global search algorithms
- Sampling uncertainty in metrics increases with newly added components
- 15

# 16 Abstract

17 SPAtial EFficiency (SPAEF) metric is one of the most thoroughly metrics in hydrologic community. In this study, our aim is to improve SPAEF by replacing the histogram match 18 19 component with other statistical indices, i.e. kurtosis and earth mover's distance, or by adding 20 a fourth or fifth component such as kurtosis and skewness. The existing spatial metrics i.e. 21 SPAtial Efficiency (SPAEF), Structural Similarity (SSIM) and Spatial Pattern Efficiency 22 Metric (SPEM) were compared with newly proposed metrics to assess their converging performance. The mesoscale Hydrologic Model (mHM) of the Moselle River is used to 23 24 simulate streamflow (Q) and actual evapotranspiration (AET). The two-source energy balance 25 (TSEB) AET during the growing season is used as monthly reference maps to calculate the 26 spatial performance of the model. The Moderate Resolution Imaging Spectroradiometer (MODIS) based Leaf area index (LAI) is utilized by the mHM via pedo-transfer functions and 27 28 multi-scale parameter regionalization approach to scale the potential ET. In addition to the real 29 monthly AET maps, we also tested these metrics using a synthetic true AET map simulated 30 with a known parameter set for a randomly selected day. The results demonstrate that the newly 31 developed four-component metric i.e. SPAtial Hybrid 4 (SPAH4) slightly outperform 32 conventional three-component metric i.e. SPAEF (3% better). However, SPAH4 significantly 33 outperforms the other existing metrics i.e. 40% better than SSIM and 50% better than SPEM. 34 We believe that other fields such as remote sensing, change detection, function space optimization and image processing can also benefit from SPAH4. 35

- 36 Keywords: mHM, model calibration, spatial pattern, SPAEF, MODIS, TSEB
- 37

### 38 1 Introduction

39 Distributed hydrologic models have a crucial role in creating digital twin of the water cycle in nature by revealing physical mechanisms and process interactions. After identifying the best 40 41 parameter set through calibration, these models are used to conduct robust numerical experiments assessing climate change impacts (Beven, 2023) or land use land cover change 42 43 impacts on model output fluxes such as runoff (Busari et al., 2021), groundwater recharge, soil moisture and actual evapotranspiration (AET). A skillful model enables decision-makers to 44 plan for and respond to water-related extremes such as hydrological droughts and floods. 45 46 Accuracy of the model results depends on the success of identifying best combination of the 47 parameters since calibration process helps us reduce discrepancies in model physics. Demirel et al. (Demirel et al., 2018) showed that using only streamflow hydrograph performance as 48 49 objective function diminishes the AET patterns simulated by the model. However, incorporating satellite based remotely sensed AET into the multi-objective calibration 50 51 framework that has already streamflow, surprisingly improves both water balance and AET 52 performance of the model. Other studies benefitted from land surface temperature (Zink et al., 2018), soil moisture (López et al., 2017; Wakigari & Leconte, 2023), AET (Avcuoğlu & 53 Demirel, 2022; Gaur et al., 2022; Odusanya et al., 2022; Sirisena et al., 2020) and groundwater 54 55 (Danapour et al., 2021; Stisen et al., 2018) in hydrologic model calibration.

56 In other words, hydrologic model calibration is essential for ensuring the validity and reliability of model predictions i.e. of most important for water management and decision-making 57 processes. However, the robustness of hydrologic model calibration heavily relies on how the 58 59 model is guided in the solution space via the performance metrics (de Boer-Euser et al., 2017; Knoben et al., 2019; Martinez-Villalobos et al., 2022; Onyutha, 2022; Schneider et al., 2022). 60 61 If the metric is too loose (tolerant) or prone to the sampling uncertainty (Clark et al., 2021), the 62 calibration process can stop quickly in the local minima while the modeler searches for the best global solution. The key point of the modelling chain is the selection of appropriate metric. 63 64 Our study focuses on development of a novel metric with least tolerance (highest 65 discrimination skill) based on benchmarking existing metrics in evaluating the similarity of two raster maps. We are particularly interested in multi-component bias-insensitive spatial 66 67 metrics for pattern comparison. Thus, bias sensitive temporal metrics used for water balance are not within the scope of this study. 68

69 The use of multi-component spatial metrics in hydrologic model calibration is an important advancement in the field of water resource management and resource allocation. The multi-70 71 component metrics provides a more nuanced evaluation of model performance compared to 72 traditional single-component metrics e.g. mean absolute error and coefficient of determination. 73 The adoption of these metrics allows for a more comprehensive understanding of the 74 hydrologic system and its spatial variability, which is critical for informed decision-making. 75 These metrics differ from single-component metrics in that they consider multiple components of the hydrological system, rather than just one component. By providing a more 76 77 comprehensive evaluation of the hydrologic system, multi-component metrics help to identify 78 areas where models can be improved. For spatial metrics, the added level of complexity 79 provided by multi-component metrics offers a more robust evaluation of model performance, 80 providing a better understanding of the spatial variability of the hydrologic system.

In recent years, remote sensing data from satellites, such as Moderate Resolution Imaging Spectroradiometer (MODIS) products, have become commonly used in hydrologic model calibration since this product provides estimates of AET from vegetation, which is a key component and major water loss in the hydrologic cycle (Becker et al., 2019; Rientjes et al., 2013). On one hand, it serves to better represent the cell-to-cell hydrological dynamics and 86 diversity in the basin also allows for a more detailed understanding of the water budget at the 87 land surface and helps to better quantify the water requirements of vegetation. On the other hand, the MODIS Leaf Area Index (LAI), product provides information about the leaf area 88 index, which is a measure of the amount of vegetation cover in an area. This information is 89 90 essential for understanding how vegetation influences the water cycle by affecting factors such as precipitation, evapotranspiration, and runoff. In this study, we use LAI to dynamically scale 91 92 the PET input to the model to improve AET performance and present a comprehensive 93 benchmarking of multi-component spatial metrics using MODIS-LAI and TSEB AET 94 products, to assess their potential for calibration (Immerzeel & Droogers, 2008).

95 There are various performance metrics in hydrology. The Nash-Sutcliffe Efficiency (NSE) and 96 Kling-Gupta Efficiency (KGE) are the most widely recognized performance metrics used in 97 evaluating and calibrating rainfall-runoff models. These two metrics have been instrumental in 98 advancing our understanding of hydrological processes and improving the performance of 99 hydrologic models (Gupta et al., 2009; Nash & Sutcliffe, 1970). They have paved the way for 100 the development of more advanced and sophisticated performance evaluation techniques. Despite the sampling uncertainty inherited in these metrics (Clark et al., 2021), NSE and KGE 101 102 continue to be widely accepted in the hydrology community due to their simplicity and 103 effectiveness in evaluating model performance. Many of the newer metrics that have been 104 introduced in recent years have been inspired by and built upon the foundation established by 105 NSE and KGE. The conventional model calibration relies on using flow-oriented temporal 106 metrics, such as the NSE and KGE. However, these metrics have a limitation as they lack 107 spatial considerations and are prone to the sampling uncertainty. This has driven the need for 108 development of intolerant spatial performance metrics which can better evaluate and improve 109 the spatial accuracy of a hydrologic model. Spatial-pattern-oriented SPAtial Efficiency (SPAEF) metric developed by Demirel et al. (Demirel et al., 2018) builds upon the strength of 110 111 KGE and incorporates new idea of distribution comparison via histogram overlap index. It is designed as a multi-component metric specifically suited for comparing spatial patterns of two 112 raster maps, with its three main data properties being co-location, variation, and distribution. 113 Although SPAEF was primarily developed for hydrologic community, it has been used in many 114 115 different disciplines such as atmospheric circulation modeling (Ahmed et al., 2019), flood risk 116 analysis (Hossain & Meng, 2020), function space optimization, fisheries (Thoya et al., 2021) and neuroscience (Yoo et al., 2020). In these studies, SPAEF has been tested and proven to be 117 118 robust and easy to interpret due to its three distinct and complementary components of 119 correlation, variance and histogram matching. Following the multi-component structure idea, we present new metrics in this study to improve SPAEF by adding fourth of fifth new 120 121 components or replacing histogram match with other components. Using this approach, we 122 aimed for reducing uncertainty in the new metric and make it sharp (discriminant) when 123 evaluating patterns on two raster maps whether they are similar or not.

124 In recent literature, there has been attempts to revise SPAEF component i.e. Spatial Pattern 125 Efficiency Metric (SPEM) (Dembélé et al., 2020). Similar to SPAEF, it has been proposed as 126 a bias-insensitive and multi-component spatial pattern-oriented metric using satellite remote 127 sensing data. Structural Similarity index (SSIM) is another pattern-oriented metric, it stands out with its spatial structure (Nilsson & Akenine-Möller, 2020; Wang et al., 2004). It was 128 129 proposed by Wang et al. (Wang et al., 2004) for image quality assessment and has been used 130 in different studies such as medical imaging, ecological restoration, and change detection in 131 the hydrological cycles and remote sensing images (Arun et al., 2021; Dougherty et al., 2020; Wiederholt et al., 2019). Knoben et al. (Knoben et al., 2019) compared NSE and KGE metrics 132 133 and argued that instead of relying directly on the KGE value, the components should be 134 analyzed in depth, even the weighting of the components. A study analyzing sampling 135 uncertainty in popular performance metrics in hydrologic modeling highlighted that the KGE can be heavily influenced by just a few data points (Clark et al., 2021). A study on the 136 hydrological model skill score compared metrics with different forms of correlation and 137 138 measures of variability, claiming the term covariance is more appropriate for evaluation 139 (Onyutha, 2022). Another recent study, based on the largest residuals, focused on reducing the largest errors, and argued that metrics should be less sensitive to errors and more sensitive to 140 141 bias (Schneider et al., 2022). The publication (Martinez-Villalobos et al., 2022) compared 142 metrics for evaluating precipitation probability distributions by comparing climate model 143 simulation data with real platform satellite data, therefore they showed the importance of 144 probability distribution functions. A study from the Netherlands (de Boer-Euser et al., 2017) 145 stated that strong components can be included in different metrics rather than considering a 146 single general metric for model comparison.

147 The existing spatial metrics aimed for the best convergence using terms such as correlation, 148 variation, histogram intersection, and root mean square error. However, kurtosis has hitherto 149 been an underrated term for spatial performance, and a four-component spatial-pattern-oriented 150 metric also does not exist for the hydrologic model calibration. We used the kurtosis ratio by 151 including it as a new component for the first time in this study in order to achieve the best 152 spatial convergence and fit. With the addition of a new component, the weighting by which the 153 components affect the value has also changed. By revealing the effect of kurtosis on spatial 154 performance, we developed a new four-component metric that does not require user input.

155 We aim to investigate the best potential to use multi-component spatial metrics in hydrological 156 model calibration, by proposing a new multi-component spatial metric that especially includes 157 the kurtosis component and benchmarking it to existing multi-component spatial metrics. The 158 primary purpose of this study is to evaluate the performance of the hydrological model using 159 multicomponent spatial metrics and to determine the potential impact on model accuracy and 160 precision. In addition, this study aims to identify the most effective combination of spatial metrics for hydrological model calibration and to develop a framework for future work in this 161 area. A large number of metrics in the literature creates confusion and difficulty for users to 162 163 choose from, so we compared metrics in this study to look for the most successful one to put a 164 stop to metric redundancy. Addressing these goals, this study aims to contribute to ongoing research efforts to improve the accuracy and reliability of hydrological models. 165

The accuracy of the analysis has been increased by comparing model predictions with real 166 platforms. It is aimed to improve the convergence between observed and simulated maps by 167 168 using two-source energy balance (TSEB) model's AET data. The MODIS-LAI data were used 169 both to correct the PET and to represent the vegetation dynamics of the Moselle basin. We utilize a spatially distributed mesoscale Hydrologic Model (mHM) with it features pedo-170 171 transfer functions for LAI data and a Multiscale Parameter Regionalization (MPR) approach to scale the potential ET (Kumar et al., 2013; Samaniego et al., 2010). We tested our framework 172 173 in three different cases to provide comprehensive outlook to the calibrations i.e. 100 iterations 174 were applied in the first case and 1000 iterations in the second case, so the effect of the number 175 of iterations was also assessed. In the third case, reproducibility was achieved by analyzing the randomly selected synthetic map. OSTRICH software (L. Shawn Matott, 2004; L.S. Matott, 176 177 2017) was used as the calibration tool and Parallel Dynamically Dimensioned Search 178 Algorithm (PDDS) was used as the calibration algorithm (Asadzadeh & Tolson, 2013). The 179 combined SPAEF value of the growing season was used as the main objective function for ET, 180 and the KGE was presented for discharge (Q) in addition. We developed multiple metrics with 181 different components and different component numbers, trying to increase the effectiveness (sharpness) of each component on convergence performance. We made an elaborated 182

183 comparison between the existing performance metrics in the literature and the newly developed 184 metrics based on SPAEF. As a result of the rigorous assessment of metrics, we identified not 185 only the superior but also new metric. The strongest aspect of this new metric is the added 186 kurtosis component.

### 187 2 Study area and data

### 188 2.1 Study area

189 The study area is the Moselle River basin, the largest part of the Rhine River basin, of which it is one of the main tributaries, characterized by diverse landforms (Figure 1). The origin of the 190 191 river from the Vosges Mountains before the interterritorial transfer from France to enter 192 Germany and Luxembourg. Furthermore, at the triangle where Germany, France and 193 Luxembourg meet, the Moselle River becomes the borderline between Germany and Luxembourg for 36 km. Also, it has a surface area of approximately 27262 km<sup>2</sup> and a length 194 195 of 545 km. Whereas, land use in the basin includes forestry, agriculture and cattle breeding in 196 the mountains and hillslopes, winegrowing on vineyards of sunny valley slopes. Moreover, the 197 altitude varies from 59 to 1326 m, with an average altitude of around 340 m (Demirel et al., 198 2013). In addition to having 26 sub-basins with surface areas varying from 102 to 3353 km<sup>2</sup>, 199 the river flow is organized by different dams, dikes, powerplants and locks such as the Trier 200 Dam, Koblenz Dam and Detzem Lock. The outlet discharge at Cochem station, located between Trier and Koblenz, varies from 14 m<sup>3</sup>/s in dry summers to a maximum of 4000 m<sup>3</sup>/s 201 202 during winter floods, with a mean discharge of around  $315 \text{ m}^3/\text{s}$  (Demirel et al., 2015).



203 204

Figure 1. DEM, land cover and AET characteristics of Mosel River basin.

An average pattern of satellite-based actual evapotranspiration for July (average of all years from 2002 to 2014) is presented to illustrate the interaction between DEM and land cover characteristics that generate the land surface flux patterns.

208 2.2 Satellite data

MODIS has a vital role in obtaining the satellite-based data used in this study, is an essential sensor aboard the Terra (EOS AM) and Aqua (EOS PM) satellites for the earth and climate measurements at a spatial resolution of approximately  $1 \text{ km} \times 1 \text{ km}$ . It provides terrestrial, atmospheric and thalassic data and a view of the entire Earth's surface for large and diverse user communities around the world. In this study, TSEB based AET is used as reference spatial

- 214 patterns (Allen et al., 1998; Norman et al., 1995). TSEB is an energy balance model using the
- energy flux principle by separating into two-layer, vegetation and soil.

216 The water limited growing season was chosen as the analysis period because it avoids climate

gradient on the AET patterns emphasizing vegetation dynamics instead of wet soil conditionsi.e. AET that is equal to the PET. All remote-sensing-based AET data were converted to long

term monthly mean data during the growing season across all years for the model calibration period (2002–2014). In what follows, three-monthly mean periods were obtained with a total

- of three-term between March and November, i.e. March-April-May (MAM), June-July-August
- 222 (JJA), and September-October-November (SON), representing AET under cloud-free
- conditions. We will attribute these AET maps as reference observations, although they are estimates from an energy balance model based on satellite observations and not pure
- 225 observations.
- Table 1. Overview of morphological and meteorological data used as input for mHM (Rakovecet al., 2016).

Variable	Description	Spatial resolution (degrees)	Source
Q (daily)	Streamflow	Point	GRDC
P (daily)	Precipitation	0.0625	E-OBS
PET (daily)	Potential evapotranspiration based on Hargreaves and Samani (Hargreaves & Samani, 1985)	0.0625	E-OBS
Tavg	Average air temperature	0.0625	E-OBS
LAI	Fully distributed 12-monthly values based on 8- day time-varying leaf area index (LAI) dataset	0.001953125	MODIS
Land cover	Forest, agriculture and urban	0.001953125	MODIS
DEM-related data	Slope, aspect, flow accumulation and direction	0.001953125	SRTM
Geology class	Two main geological formations	0.001953125	ESD UFZ – Leipzig (Rakovec et al., 2016)
Soil class	Fully distributed soil texture data	0.001953125	HWSD

228 GRDC - Global Runoff Data Centre, E-OBS - The gridded observational dataset from Copernicus, MODIS - Moderate

229 Resolution Imaging Spectroradiometer, SRTM – Shuttle Radar Topography Mission, ESD – European Soil Database, HWSD –

230 Harmonized World Soil Database

### 231 **3 Hydrological model**

232 This research utilizes the mesoscale Hydrologic Model (mHM) v.5.11.2 (Samaniego et al., 233 2021) which is a grid-based spatially distributed model it features pedo-transfer functions and 234 MPR (Kumar et al., 2013; Samaniego et al., 2010; Thober et al., 2019). Another feature of 235 mHM is the use of leaf area index (LAI) data not only for calculating interception loss but also 236 for dynamically scaling PET (Demirel et al., 2018). With these unique features, it is more 237 flexible than other existing hydrologic models in line with the purpose of this study. The model 238 features 69 adjustable global parameters that can be optimized during the calibration process 239 (Demirel et al., 2018). The model works on the basis of water balance rather than energy 240 balance and provides various physically meaningful spatial outputs, fluxes and states as 241 simulating major elements of the hydrologic processes, i.e. soil moisture dynamics, 242 interception, infiltration, evapotranspiration, snow accumulation and melting, groundwater storage, seepage, surface runoff and others. 243

244 The basic data for the running mHM can be classified into meteorological data, morphological 245 data, land cover data and gauge streamflow data. Table 1 shows a summary of the data used in 246 mHM setup provided by Rakovec et al. (Rakovec et al., 2016). As seen in the table, mHM can 247 handle different spatial resolutions of meteorological data and morphological data since it has internal upscaling and downscaling subroutines. At this point, the Multi-Scale Parameter 248 249 Regionalization technique comes into play and enables user to map calibrated parameters to 250 the simulated grids with pedo-transfer functions. This approach prevents uniform parameter 251 fields and protects sub-grid heterogeneity of the fluxes. In other models, every parameter gets the same value in the entire sub-basin or in each hydrologic response units resulting in uniform 252 253 flux results for the same domain.

254 The meteorological model inputs are precipitation, average air temperature and potential 255 evapotranspiration (PET). In our study, PET was direct input to the mHM and estimated outside 256 with Hargreaves-Samani (Hargreaves & Samani, 1985) method using additional temperature 257 data. All meteorological data are obtained from E-OBS at daily resolution, originally at 10-20 258 km. The morphological variables are digital elevation model (DEM), soil maps with textural 259 features, geological maps including specific yield, permeability and aquifer thickness. In 260 addition to characterizing the morphology of the basin, DEM masks the grid cells with the 261 basin boundaries to eliminate no-data parts. All morphological data are prepared at 262 0.001953125 degrees ( $\sim$ 200 m  $\times$  200 m) scale. The model hydrology is evaluated at 0.015625 degrees (~2x2 km) spatial resolution and daily time step. Lastly, monthly leaf area index (LAI) 263 264 maps are used to represent the vegetation dynamics for both interception calculation and PET 265 correction for the entire period (2002–2014). Four years of model warm-up period (1998–2001) 266 is used. Observed daily streamflow (Q) data at Cochem (station #6336050), provided by the Global Runoff Data Centre (GRDC), Koblenz (Germany), is used to calibrate water balance in 267 the basin. 268

### 269 4 Methods

270 In this study, we tested nine different spatial metrics i.e. two of them are existing metrics, and 271 seven of them are newly developed based on SPAEF (Table 2). To evaluate the effect of 272 number of iterations, calibrations were pursued with either 100 or 1000 maximum iterations. 273 Besides, synthetically created AET maps using mHM and a pre-defined parameter set are 274 utilized to mimic a "hide and seek" case. This is crucial to test the guidance performance of the 275 metrics in the multi-dimensional solution space to find the hided (perfect) solution within 1000 276 iterations since search algorithms, i.e. ParaPADDS algorithm herein, require a metric to 277 evaluate model results at every iteration.

### 278 4.1 Objective Functions

Multi-component structure of our metrics was inspired by the Kling–Gupta efficiency (Gupta et al., 2009). KGE is one of the most used metrics in the hydrologic modelling to evaluate streamflow performance. As shown in Eq.(1), it has three components, i.e., correlation, variability and bias.

$$KGE = 1 - \sqrt{\left(\alpha_Q - 1\right)^2 + \left(\beta_Q - 1\right)^2 + \left(\gamma_Q - 1\right)^2}$$

$$\alpha_Q = \rho(o, s), \beta_Q = \frac{\sigma_S}{\sigma_0} \text{ and } \gamma_Q = \frac{\mu_S}{\mu_0}$$
(1)

where  $\alpha_Q$  is the Pearson correlation coefficient between the observed (o) and the simulated (s) discharge time series,  $\beta_Q$  is the relative variability based on the ratio of standard deviation in simulated and observed values and  $\gamma_Q$  is the bias fraction which is normalized by the standard deviation of the observed data.

Table 2 shows the summary of SPAEF based metrics. For brevity, we used Eq. (2) as formula template i.e. a generic formulation type that encompasses in the number and content of components. The excessed style in Eq (2) includes all metrics form with various components.

METRIC = 
$$1 - \sqrt{(\alpha - 1)^2 + (\beta - 1)^2 + (\gamma - 1)^2 + (\kappa - 1)^2 + (\delta - 1)^2}$$
 (2)

Matuin	Components								
	α	β	γ	к	δ				
SPAtial Efficiency (SPAEF)	$\rho(o,s)$	$\frac{\sigma_o}{\mu_o}/\frac{\sigma_s}{\mu_s}$	$\frac{\sum_{j=1}^{n} \min(K_j, L_j)}{\sum_{j=1}^{n} K_j}, n=100 \text{ fixed}$	none	none	Eq. (3)			
SPAtial EFficiency Prime (SPAEF')	same as SPAEF	same as SPAEF	same as SPAEF except for dynamic <i>n</i> i.e. number of bins $n = floor\{\sqrt{length(o)}\}$	none	none	Eq. (4)			
SPAtial Count Density Efficiency (SPACD)	same as SPAEF	same as SPAEF	$\frac{\sum_{j=1}^{n} \min(K_j, L_j)}{\sum_{j=1}^{n} K_j} \left( v_n = c_n / w_n \right)$	none	none	Eq. (5)			
SPAtial Hybrid 4 Efficiency (SPAH4)	same as SPAEF	same as SPAEF	same as SPAEF'	Kurt(s) Kurt(o)	none	Eq. (6)			
SPAtial Kurtosis Efficiency (SPAK)	same as SPAEF	same as SPAEF	none	same as SPAH4	none	Eq. (7)			
SPAtial Hybrid 5 Efficiency (SPAH5)	same as SPAEF	same as SPAEF	same as SPAEF'	same as SPAH4	Skew(s) Skew(o)	Eq. (8)			
SPAtialHistogramEqualizationEfficiency(SPAHE)Efficiency	same as SPAEF	same as SPAEF	$\frac{\sum_{j=1}^{n} min(K_j, L_j)}{\sum_{j=1}^{n} K_j}$	none	none	Eq. (9)			
SPAtial Movers' Distance Efficiency ( <b>SPAMD</b> )	same as SPAEF	same as SPAEF	$\frac{\sum_{i=1}^{K} \sum_{i=1}^{L} f_{i,j} d_{i,j}}{\sum_{i=1}^{K} \sum_{i=1}^{L} f_{i,j}}$	none	none	Eq. (10)			
Spatial Pattern Efficiency Metric ( <b>SPEM</b> )	$1 - \frac{6\sum_1^n d^2}{n(n^2 - 1)}$	same as SPAEF	$1 - E_{RMS}(Z_{X_s}, Z_{X_o})$	none	none	Eq. (11)			

 Table 2. SPAEF based metrics used as objective functions.

291

293 SPAEF is the seed of our newly proposed metrics as our aim is to sharpen SPAEF. In other 294 words, we intend to improve its discriminating power while judging whether two maps are 295 similar or not. SPAEF uses a multi-component structure of the KGE metric. In Eq. (3),  $\alpha$  is the 296 Pearson correlation coefficient between the observed (o) and simulated (s) pattern,  $\beta$  is the 297 fraction of the coefficient of variation representing spatial variability and  $\gamma$  is the histogram 298 intersection, which based on z-scores, for the given histogram K of the observed pattern and 299 the histogram L of the simulated pattern, each containing n bins (Swain & Ballard, 1991). The 300 SPAEF can have a value between  $-\infty$  and 1, where a value closer to 1 indicates highest spatial 301 similarity between the observations and model simulations (Koch et al., 2018).

As a result of various adjustments and improvements made in the SPAEF components, new
 metrics were proposed and tested i.e. SPAEF', SPACD, SPAH4, SPAK, SPAH5, SPAMD, and
 SPAHE. We included two popular metrics, SPEM and SSIM into benchmark.

First improvement in SPAEF is changing user defined the number of bins to an automated n based on the number of elements (grids) in the raster map (see Eq (3)). We introduced a simple approach i.e. the square root of the length of the observed data as n = $floor{\sqrt{length(o)}}$  although there are different methods for the same purpose (Freedman & Diaconis, 1981; Scott, 1979; Sturges, 1926). This slightly new version of the SPAEF is presented as SPAEF-Prime (SPAEF') as shown in Eq (4). Unlike the standard version, the SPAEF' does not require any user-defined inputs now.

Eq (5) shows Spatial Count Density Efficiency (SPACD) which has a different type of normalization based on count density approach in the calculation of the histogram intersection component. While the first two components remain constant as in SPAEF' the calculation of n in the gamma component has changed. This approach uses count or frequency scaled by the width of the bin  $v_n = c_n / w_n$ ,  $v_n$  is the bin value,  $c_n$  is the number of elements in the bin and  $w_n$  is the width of the bin, respectively.

Eq (6) shows SPAtial Hybrid 4 Efficiency (SPAH4) which is a four-component metric obtained 318 319 by adding kurtosis i.e. a fundamental statistical property of distributions to the SPAEF' metric. 320 Kurtosis can be defined as a measure of how prone a distribution is to outliers (Pearson, 1905). 321 SPAH4 offers a more accurate perspective by questioning not only the match of the histograms but also the extreme values and spread in the data. The 4<sup>th</sup> component is symbolized by the 322 323 expression Kurt and  $\kappa$  is the ratio of the kurtosis coefficients of the simulated (s) and observed 324 (o) data. Eq. (7) shows SPAtial Kurtosis Efficiency (SPAK) which is a three-component metric 325 replacing the histogram intersection component in the SPAEF metric with the kurtosis 326 coefficient component. Thus, it dominates the metric on its affinity for discrete values without 327 questioning histogram intersection.  $\alpha$  and  $\beta$  were introduced and explained in previous metrics, 328 also  $\kappa$  is declared in Eq. (7) as ratio of kurtosis coefficient. This metric can be characterized as 329 a mixture of SPAH4 and SPAEF metrics. Eq. (8) shows SPAtial Hybrid 5 Efficiency (SPAH5) which is a five-component metric adding skewness to the SPAH4 metric. Skewness can be 330 331 defined as a measure of the asymmetry of the data around the sample mean.

Eq. (9) shows SPAtial Histogram Equalization Efficiency (SPAHE) that is very similar to SPAEF with additional step before histogram match calculation "histogram equalization" approach. This approach is a computer image processing technique used to improve contrast in raster data. Its quantitative logic is based on the grayscale transformation (*T*) to minimize  $|c_1(T(k)) - c_0(k)|, c_0$  is the cumulative histogram of the input data, and  $c_1$  is the cumulative sum of target histogram for all intensities *k*. Histogram equalization is a specific case of the histogram remapping methods. It is an image processing technique used to advance contrast in images which spatial patterns for this study. It achieves this by efficaciously sprawling out the
 most frequent intensity values, i.e. expanding the intensity range of the image (Efford, 2000).

Eq. (10) shows SPAtial Efficiency Movers' Distance (SPAMD) is another SPAEF-oriented 341 342 multi-variate metric which measures the quantitative closeness of two pattern set by 343 considering the Earth Movers' Distance of their histograms (Rubner et al., 1998). The aim of 344 EMD approach is minimization of overall transfer cost in the conversion one histograms to another. In Eq (10),  $f_{i,i}$  is flow cost of transfer ith term of histogram K of observed map to jth 345 histogram L simulated map at distance  $d_{i,i}$ . EMD is the ratio of work done through the total 346 optimal flow and the total flow. The value of EMD is zero indicates the perfect consistency 347 348 between two histograms.

349 Eq. (11) shows Spatial Pattern Efficiency Metric (SPEM), a metric inspired by KGE and 350 SPAEF, is one of the existing metrics included in our analysis (Dembélé et al., 2020). It forces 351 the z-scores of simulated variables and observed variables to be equal (i.e., minimizing their 352 ERMS) corresponds to matching their grid cell locations (i.e., spatial patterns). SPEM 353 considers a modeled variable (Xmod) and an observed variable (Xobs) of n elements, it is 354 defined as Eq. (11); where rs is the Spearman rank-order correlation coefficient with d the difference between the ranks of Xmod and Xobs.  $\gamma$  is the variability ratio that assesses the 355 356 similarity in the dispersion of the probability distributions of Xmod and Xobs, with  $\mu$  and  $\sigma$ representing the mean and the standard deviation, respectively, and  $\alpha$  the spatial location 357 358 matching term calculated as the root-mean-square error (ERMS) of the standardized values (z-359 scores, ZX) of Xmod and Xobs (Dembélé et al., 2020). The formula for d can be written as  $d = diff(rank(X_s), rank(X_o))$ . SPEM ranges from  $-\infty$  to 1, which is its optimal value. 360

Lastly, Eq. (12) shows Structural Similarity index (SSIM) (Wang et al., 2004). An image 361 quality metric SSIM to evaluate degradation grade caused by visual data processing. This 362 363 method considers pattern similarity as it detects changes in the variation of structural 364 information between the two images. The algorithm formulates perception sensibility to visual changes based on the distortion luminance, contrast and structure information. By combining 365 three components, similarity can be characterized with overall unit metric in terms of statistical 366 properties of simulated and observed data such as mean  $\mu$ , standard deviation  $\sigma$  and covariance 367  $cov_{o,s}$ , as shown in Eq. (12).  $c_1$ ,  $c_2$  are constants that stabiles functions when the dominator 368 terms are close to zero. The SSIM is a fully referenced objective quality metric that gives values 369 370 in the range [0,1] relative to the structural relationship between the two images.

$$SSIM = \frac{(2\mu_o\mu_s + c_1)(2cov_{o,s} + c_2)}{(\mu_o^2 + \mu_s^2 + c_1)(\sigma_o^2 + \sigma_s^2 + c_2)}$$
(12)

371

All nine spatial metrics were calculated separately as long term (2002-2014) monthly average of AET data for three periods covering the growing season and combined as in Eq (13) to minimize the total error, representing objective function (OF). These periods are symbolized as March-April-May (MAM), June-July-August (JJA), and September-October-November (SON).

$$Minimize \left[ (1 - METRIC_{MAM})^2 + (1 - METRIC_{JJA})^2 + (1 - METRIC_{SON})^2 \right]$$
(13)

377 It should be noted that although we tested other metrics and approaches, we only reported nine 378 selected metrics in this study. For instance, we used harmonic mean or geometric mean instead 379 of the arithmetic mean in the second component of SPAEF. In another attempt, we replaced 380 the skewness coefficient ratio with different L-moments. We also used Hausdorff distance 381 (Hausdorff, 1914) and Fréchet distance (Fréchet, 1906) as third component in SPAEF. Even we used the product of components i.e. multiplied them instead of adding them. However, all 382 383 these attempts did not reveal better results than those reported in this study. Therefore, for 384 brevity we reported the ranking of only these nine metrics above. In this calibration study, we fine-tuned only 20 parameters of daily mHM for the Mosel Basin using the popular global 385 386 search algorithm Pareto-Archived Dynamically Dimensioned Search (ParaPADDS) algorithm (Asadzadeh & Tolson, 2013) using 750 maximum iteration and 3 parallel cores. The 20 387 388 parameters out of 69 mHM parameters are selected based on a sensitivity analysis done in our 389 previous study. Note that ParaPADDS is the multi-objective version of the Dynamically 390 Dimension Search algorithm (Tolson & Shoemaker, 2007) available in OSTRICH 391 Optimization Software Toolkit (L.S. Matott, 2017).

### 392 **5 Results**

393 In this study, six novel metrics are proposed and compared with existing SPAEF, SPEM and 394 SSIM metrics in pattern analysis of distributed hydrologic model simulations. The new metrics 395 can be called as "the sisters of SPAEF" as they have emerged from the well-established SPAEF 396 with additional unique statistical features such as automated number of bins, kurtosis and 397 skewness included in their structure. We ranked the nine metrics based on their effectiveness 398 in distinguishing between two raster maps during distributed model calibration with MODIS-399 LAI and TSEB AET for a period of 13 years from 2002 to 2014. Pre-selected 20 mHM 400 parameters are included in the following three different pattern-only calibration cases: (1) 100 401 iterations with satellite data, (2) 1000 iterations with satellite data, and (3) 1000 iterations with 402 synthetic maps. Synthetic map represents a map simulated with a known mHM parameter set for a randomly selected day that is used as the target in parameter optimization (calibration) 403 process. The use of this synthetic scenario is planned to ensure the reproducibility of the 404 405 analysis and to have a fully controlled numerical experiment. Obviously, long term monthly 406 averaging was done only with real satellite data to form robust seasonal pattern maps i.e. target 407 in the calibration.

408 Although water balance metrics, i.e. temporal metrics, are not included in the calibration, KGE 409 values are calculated to evaluate the model simulations together with standard SPAEF in Table 410 3. Streamflow simulation performance was calculated for the calibration period (2002-2014), 411 using the KGE metric between the observed gauge streamflow and simulated streamflow from 412 the model. This is done only for case 1 (TSEB 100 runs) and 2 (TSEB 1000runs) i.e. real satellite data are used in the pattern-based optimization. It is interesting to note that some of 413 414 the pattern metrics help to improve the bias in water balance as well. The three OF columns in 415 this table show lowest (best) values of each metrics reached using Eq. (13). This is particularly 416 important to show the skill of the nine metrics in converging to zero i.e. certainly exists in the synthetic case (3). It should be noted that the metrics are ranked based on the standard SPAEF 417 418 values. Closer inspection of the Table 3 shows that TSEB 1000 iterations significantly 419 improves the SPAH4 performance from 0.608 to 0.688 (SPAEF value) as compared to the TSEB 100 iterations. The reduction in OF is even more remarkable since the error in SPAH4 420 421 was halved from 0.70 to 0.35 when iterations are increased to 1000. It is clear from this table 422 that SPAHE and SPAH5 are the worst performing two metrics among all three cases.

423 Comparing the two results (100 runs vs 1000 runs) it can be seen that all metrics are improved 424 with the increased number of iterations showing the importance of the selecting appropriate 425 number of the iterations for the search algorithm. However, if enough freedom is not given to 426 the optimizer, it may fail to find the global optimum point in the solution space. Combining

427 kurtosis with skewness in the same metric (SPAH5) did not produce a discriminative metric.

428 This result is somewhat counterintuitive as we expect more constrain would yield improved

429 performance. What is striking about the values in this table is histogram equalization step did

430 not help to improve the pattern results and discriminative power of the metric.

431 **Table 3.** Calibration results of the three cases. Note that metrics are ranked based on 1000 run
432 - SPAEF values (4<sup>th</sup> numeric column).

Metrics	TSEB 100 runs			TSEF	8 1000 run	S	SYNTHETIC MAP 1000 runs		
	SPAEF	KGE	OF	SPAEF	KGE	OF	SPAEF	OF	
SPAH4	0.608	0.78	0.70	0.688	0.77	0.35	0.948	0.05	
SPACD	0.619	0.26	0.40	0.673	0.74	0.27	0.939	0.04	
SPAEF'	0.585	0.36	0.52	0.671	0.52	0.33	0.949	0.05	
SPAK	0.558	0.89	0.39	0.638	0.87	0.25	0.906	0.01	
SPAMD	0.614	0.07	0.29	0.625	0.66	0.21	0.859	0.02	
SSIM	0.557	0.21	0.19	0.491	0.41	0.15	0.948	0.00	
SPEM	0.609	0.33	1,71	0.460	0.61	1,46	0.941	0.05	
SPAHE	0.492	0.70	0.25	0.376	0.65	0.21	0.758	0.04	
SPAH5	-0.519	0.61	8,15	0.211	0.53	2,07	0.953	0.05	

### 433

434 What stands out in the table is that SSIM seems to be the most tolerant metric reaching lowest 435 OF values which corresponds to the poor SPAEF performance in all three cases. In case 3, in 436 particular, the search algorithm could converge nearly to zero SSIM but the evaluation of the 437 maps with SPAEF revealed that it is only a match around 0.95 SPAEF and not very close to 1 SPAEF i.e. perfect pattern match. In other words, minimizing SSIM in Eq (13) nearly to zero 438 439 after calibration doesn't guarantee a perfect pattern match in terms of SPAEF currency 440 (metric). Based on the results of case 1 and 2, SPAH4 and SPAK are the most successful spatial metrics for water balance. Obviously, SSIM and SPAMD have the worst KGE performance in 441 442 case 1 and 2. Note that KGE is not calculated for the synthetic case 3. Interestingly, the 443 minimization of SPEM and SPAH5 metrics via Eq (13) after optimization resulted in poor 444 values above 1 both in case 1 and 2.

445 Figure 2 shows the reference AET maps and simulated AET maps from the mHM with 446 calibrated parameters after 100 iterations (case 1). The reference three maps are given in both 447 columns for ease of comparison. The order of the metrics is in accordance with the performance 448 ranking in Table 3 and also, the ranking is provided (e.g. #1, #2 etc.) to help to the reader. The 449 combined SPAEF values of three periods (MAM, JJA and SON) are presented in brackets 450 underneath the metric name. To use a single legend, the maps are normalized with their mean. 451 The resultant maps from SPACD and SPAMD (second row in Figure 2) are slightly better than other rows as visually more similar to the reference maps (first row in Figure 2). Closer 452 453 inspection of the maps shows that the high contrast between west and south of the basin in SON period is well-captured by most of the metrics except for the SPAH5 (row 6, rank #9). 454



458 Figure 2. Long term average three-monthly TSEB reference maps versus mHM simulated 459 maps using MODIS-LAI and best-balanced Pareto solution parameter set from 100 run case.

460 Figure 3 shows the reference AET maps and simulated AET maps from the mHM with 461 calibrated parameters after 1000 iterations (case 2). It is consistent with Figure 2 that the 462 simulated AET maps by the model parameter sets optimized with SPAH4 and SPACD metrics 463 are most close to the reference maps. Similarly, the poor AET performance of SPAH5 maps is 464 apparent from the maps in the last row of the figure. Map illustration of each period reveals 465 that the combined metric value (OF) can hinder individual map performance. For instance, the SON map of the SPAHE metric in Figure 3 shows that the model better converges to the 466 remotely sensed reference map when optimized with SPAHE whereas the MAM and JJA maps 467 468 show that the model could not reproduce the AET maps of these periods as successful as with 469 the other metrics.



472 Figure 3. Long term average three-monthly TSEB reference maps versus mHM simulated 473 maps using best-balanced Pareto solution parameter set from 1000 run case.

474 The entire calibration development process, the model improvements from beginning to end 475 and the optimum points are depicted with scatter diagrams in Figure 4. It shows the relationship 476 between the value and iteration based on the ParaPADDS search algorithm, more specifically, 477 the objective function value achieved for each iteration step of the calibration process. While 478 the OF results in Table 3 are obtained at the end of the iteration step sequence, some consistent 479 metrics may reach this best value earlier. SPAH4 reached its best OF value at 0.70 and 0.35 in 480 approximately quarter steps for 100 and 1000 runs, respectively. Similarly, SPACD, SPAEF' 481 and SPAMD are also fast-improving metrics. Since the synthetic case was based on a virtually 482 generated daily map, it took longer for the metrics to find the points where their improvement 483 became linear, nearly a third. It is surprising to see that SPAH5 and SPEM are consistent early 484 maturing metrics despite their poor spatial performance.







Figure 4. Scatter plots of the calibration processes, the OF value-iteration relationship of the 488 489 PDSS search algorithm. First and second column sub-plots are the same figures except for 490 different extent.



492 Figure 5. Monthly average hydrograph of all years in the calibration period (2002–2014) to
 493 demonstrate the flow simulation performances of nine different metrics.

494 Figure 5 compares in-situ observed hydrographs and simulated hydrographs constrained by 495 metrics. SPAH4 and SPAK performed better in each case, predicting the most similar 496 discharges to the observed Cochem outflows. Otherwise, the SPAHE metric standout for the 497 100 runs and the SPACD metric for the 1000 runs, as pointed out by the KGE column in Table 498 3. The simulations show better hydrograph fitting during the growing season, especially during 499 the summer months, also the hydrograph line breakpoints, peaks and valleys are coherently 500 followed. Thus, the overall trend and characteristics of the streamflow were successfully 501 analyzed and represented. Also, a positive correlation was found between increasing iteration 502 and hydrograph fit. As the number of iterations increases, the hydrograph lines become closer to the observed lines and the overall consolidation of the hydrographs provides better results. 503 504 The narrow range of hydrographs in Figure 5 shows that the developed new metrics can be 505 used not only for the spatial pattern performance simulating the AET but also for the temporal 506 streamflow performance simulating the discharges.

507 Overall, the results indicate that the newly developed SPAH4 and SPACD are the best 508 performing metrics for all calibration scenarios, particularly in the non-synthetic TSEB cases. 509 The competitive performance of the SPAMD metric that follows them should not be ignored. 510 Briefly, the four-component spatial performance metric SPAH4 stands out especially with its 511 versatile evaluation and robust performance, indicated with bold text in Table 3. Although the 512 modeler can use the SPAH4 and SPACD metrics in the long and short runs, respectively, both 513 offer close values for the decision makers. We can see that the only negative output is experienced in the TSEB 100 runs i.e. SPAH5 It should not be overlooked that SPAH5 is a 514 515 prominent metric for synthetic scenarios. Interestingly, there is a significant positive correlation 516 between the KGE and the metrics containing the kurtosis statistic.

### 517 6 Discussion

518 This study sets out to assess the importance and comparison of spatial metrics in distributed 519 model calibration. Previous studies have noted that spatial metrics are closer to the reference model than time series metrics in model optimizations (Demirel et al., 2018). One of the first 520 521 objectives of the study is to select the appropriate spatial performance metric that plays an 522 active role in simulating inadequate spatial AET models similar to satellite-based reference 523 models. SPAEF has been the inspiration for this study with its innovations in spatial model parameterization and spatial performance metric selection. These innovations have raised new 524 525 questions in the pattern comparison used in model optimization. Numerous imperfect models 526 are produced during these optimizations, due to limitations in the chosen objective function. 527 To overcome these limitations and to obtain a more physically meaningful and empowered 528 metric, we have developed new metrics that include statistical and analytical approaches. 529 Thanks to this meta-analysis, while suggesting the most successful metric for users, different 530 objective functions that can be used for various purposes can also be seen as an opportunity. 531 While searching for new solutions for a more robust spatial performance metric, we derived 532 metrics that emphasize spatiality in a more comprehensive way by increasing the number of 533 components of SPAEF and changing the content of the components. For the three cases, 534 significant findings that are both different from each other and support each other have been 535 identified. The TSEB 100 and 1000 run cases in model calibration served the purpose of 536 evaluating metric performances in short and long runs, thus providing a flexible and versatile 537 assessment that allows the progress of the model calibration performed by the metrics to be 538 monitored and the decision maker to choose metrics according to their preferences.

539 TSEB 100 runs, which we tested by focusing on the performance of spatial metrics in short 540 runs, SPACD and SPAMD demonstrated better results on the SPAEF basis compared to other 541 metrics. Notably, SPAK and SPAH4 including the kurtosis coefficient ratio component, 542 yielded the best KGE values even at iterations close to the beginning. TSEB 1000 runs which 543 we tested by focusing on its performance in long runs, resulted in more decisive outcomes with 544 no negative values for any criteria. SPAH4 emerged as the top-performing metric in this case, 545 followed by SPACD. The competition between these metrics was notable. In the uncertainty 546 analysis, SPAH4 has an acceptable sampling error although it has the extra component. (Table A1). Like the TSEB 100 runs, SPAK and SPAH4 exhibited the highest KGE values. This 547 548 indicated consistency was strong evidence for important findings and suggests that the 549 descriptive statistical kurtosis ratio component has a considerable positive effect on the discharge simulation. Due to the tendency of the SPAH4 metric including kurtosis for flow 550 551 prediction, it worked as a metric that focused on both spatial and flow performance, although 552 the analysis was performed with a single spatial performance-oriented objective function. It 553 sheds light on the analysis in detecting the presence of outliers potential also differences in the tail and crests, controlling data integrity, understanding data distribution, reliability of the 554 555 statistical analysis and improving the metric performance from a statistical perspective. Thus, by investigating and questioning the effect of outliers on spatial performance, the harmony and 556 557 differences between them are also included in the model. Now that these outliers are introduced to the model, the histogram intercept component is also supported, the margin of error is 558 559 reduced and a more exact match is made.

560 In the synthetic scenario, the metric SPAH5 which incorporates skewness characteristics, 561 yielded the best SPAEF value. SPACD and SPAH4 also demonstrated successful outcomes in 562 this scenario. The kurtosis information we use in the SPAH4 metric expresses how often 563 outliers occur, while the skewness information we use as the fifth component in the SPAH5 564 metric gives information about the direction of the outliers. Our purpose in including the 565 skewness component is to question the likelihood of events in the probability distribution, and 566 especially to consider extreme distribution. Various datasets have different characteristics, since the differences specific to this dataset represent important concepts in the calibration 567 568 model, many principles are referred to using the skewness information, from the algorithm of the model to the physics-based hydrology information. Thus, we enabled a more 569 570 comprehensive and more specific analysis for models consisting of diverse data. Our finding 571 of the importance of these statistical measures in understanding the data is supported by the 572 study by Cain et al., processing skewness and kurtosis information on distributions collected 573 from the authors of the published articles (Cain et al., 2017). In addition, it is possible to derive 574 a positive interpretation from a negative finding in meta-analyses as in this study. Since the 575 only difference between the metrics with the best and the worst performance in TSEB runs, 576 namely SPAH4 and SPAH5, is the skewness ratio component, it can be concluded that 577 skewness is a component that negatively affects the spatial metrics used in pattern comparison. 578 It should be noted though that skewness information is an outstanding component for synthetic 579 cases.

580 In TSEB 100 runs scenario, the spatial performance tussle results of the metrics show that the 581 newly proposed metric i.e. SPACD outperforms the conventional three-component metric 582 SPAEF (5.76% better) on the other hand 11.11% better than SSIM and 1.66% better than 583 SPEM. In TSEB 1000 runs results demonstrate that the newly developed four-component 584 metric i.e. SPAtial Hybrid 4 (SPAH4) slightly outperform SPAEF (2.62% better). However, 585 SPAH4 significantly outperforms the other existing metrics i.e. 40.22% better than SSIM and 586 49.53% better than SPEM.

### 587 7 Conclusion

588 In this study, we thoroughly assessed common existing metrics and new spatial pattern-oriented performance metrics that we developed based on SPAEF. For the consistency and reliability 589 590 of the results, the Mosel Basin with high data quality was selected and the physics-based fully 591 distributed mHM model was established for this basin. In these three different scenarios, we 592 performed analyses with various (low-high) iterations for actual evapotranspiration maps (TSEB AET) and synthetic maps. The most popular metrics (SPAEF, SSIM and SPEM) were 593 594 compared with new metrics (SPAH4, SPACD, etc.) to measure the convergence of the mHM 595 model to long-term monthly AET maps observed during parameter calibration. The usage of 596 this synthetic scenario is important to ensure the reproducibility of the experiments and to give 597 us full control over the calibration process. Based on our findings we can draw the following 598 conclusions.

The inclusion of kurtosis ratio coefficient in the spatial pattern-oriented metrics demonstrates
 that metric performance is improved, so it has a positive impact on the spatially objective
 functions. Also shows a positive effect on streamflow prediction, it successfully calibrates the
 KGE metric even in very short runs. Furthermore, while using the skewness ratio coefficient
 gave unsuccessful results for TSEB AET maps, the kurtosis information of the distribution was
 more prominent in the pattern performance of the models. However, the SPAH5 performs the
 best among the close results and is presented as a strong hypothesis for the synthetic cases.

- The metric with the best performance in the short runs was SPACD, which normalizes the
distribution according to density. The excellent consistency between histograms, which is the
main component of the Earth mover's distance metric, has a positive effect on making this
metric a sharp metric with little tolerance, making SPAMD the second-best metric.

- 610 The best-performing metric on long runs was SPAH4, a four-component spatial performance
- 611 metric that includes the kurtosis of the distribution. It was followed by the SPACD metric,
- 612 which proved its consistent performance. Thus, the decision maker is presented with a flexible
- 613 and wide working area.
- Considering all the experimental results, the most successful and robust metric in all three
   scenarios is our newly developed spatial pattern-oriented SPAH4, which outperforms the
   existing metrics in the literature by up to fifty per cent.
- 617 In future studies, it would significantly enhance the depth and quality of the analysis to increase
- 618 the number of iterations. In fact, convergence in hydrological models is closely related to the
- 619 number of parameters and the freedom of the appropriate iteration chosen. Future work may
- benefit from exploring untested statistical terms to add a new perspective. We expect that these
- 621 newly developed metrics, especially SPAH4, will be used not only in hydrology but also in
- 622 other fields including remote sensing, image processing and object detection.

623 **Appendix A:** Results of the jackknife and bootstrap based sampling uncertainty analysis. Clark et al (2021) showed that the two most popular metrics in hydrology, i.e. NSE and KGE, are 624 625 vulnerable to sampling uncertainty since the differences between observed and simulated 626 streamflow values at random time steps in time series which can have significant effects on the results (Knoben & Spieler, 2022). From this study, we are inspired to assess the sampling 627 uncertainty in ten metrics using the gumboot R package (Clark et al., 2021) which uses a 628 629 jackknife-after-bootstrap method of Efron (1992) to estimate standard errors (SEJaB) shown 630 in Table A1.

GOF_stat	seJack	seBoot	p05	p50	թ95	score	biasJack	biasBoot	seJab
SSIM	0.0103	0.0099	0.6144	0.6311	0.6457	0.6307	-0.0002	0.0000	0.0091
SPAHE	0.0568	0.0119	0.7717	0.7917	0.8107	0.7783	0.1496	0.0131	0.0112
SPAMD	0.0114	0.0108	0.6785	0.6972	0.7137	0.6966	0.0006	0.0001	0.0115
SPEM	0.0180	0.0175	0.2739	0.3041	0.3309	0.3034	-0.0006	-0.0003	0.0146
SPAEF	0.0133	0.0128	0.6489	0.6711	0.6917	0.6727	0.0017	-0.0021	0.0148
SPAEF'	0.0133	0.0127	0.6489	0.6711	0.6917	0.6727	0.0017	-0.0021	0.0152
SPAK	0.0302	0.0288	0.5719	0.6226	0.6661	0.6207	-0.0004	0.0007	0.0295
SPAH4	0.0302	0.0298	0.5484	0.5999	0.6459	0.6000	0.0011	-0.0012	0.0313
SPACD	0.0234	0.0248	0.6056	0.6571	0.6851	0.6670	-0.0219	-0.0142	0.0603
SPAH5	0.1685	0.2077	-0.3594	0.0373	0.2636	0.0427	-0.0388	-0.0401	0.3382

631 **Table A1.** Sampling uncertainty of the metrics i.e. ranked based on the seJab column.

632



Figure A1. eCDF plot of daily discharge for all years in the calibration period (2002-2014) to
 visualize the distribution of the data and identify statistical patterns.

637 Figure A1 visualizes the empirical cumulative distribution function (eCDF) plot for the 638 observed and simulated data, which shows how the probability of a given discharge value 639 occurring varies over the range of discharge values. In this context, the percentage of observed 640 discharges less than nearly 500 is 80% and less than 200 is 50% for both the TSEB 100 and 641 1000 runs. Furthermore, the slope of the curve at any point represents the density function of 642 the discharge values at that point, and the intervals where the curve steepens contain values 643 close to the mean value. Hence, it can be concluded that the overall average discharge value of the steepening intervals of the flow data resulting from the simulation of the metrics is roughly 644 645  $300 \text{ m}^3$ /s. The mean observed outflow of Cochem station is around  $315 \text{ m}^3$ /s supports this 646 outcome. In both cases, SPAK and SPAH4 illustrated a high level of matching in terms of the 647 fit of the curves generated by the observed data (OBS) and the metrics, with the least difference 648 between the distributions.





Figure A2. Monthly average hydrograph of the last two years in the calibration period (2013–2014)

### 654 Acknowledgements, Samples, and Data

655 Data Availability Statement: Discharge data is provided by GRDC data portal (https://portal.grdc.bafg.de/) in Koblenz, Germany. MODIS MOD16A2 v061 product was 656 retrieved from https://doi.org/10.5067/MODIS/MOD16A2.061. SRTM DEM data was 657 retrieved from https://www.earthdata.nasa.gov. The source code of the mHM is publicly 658 659 available at https://doi.org/10.5281/zenodo.4575390. The source code of the SPAEF metric is publicly available at https://doi.org/10.5281/zenodo.5861253. The model calibration software 660 Ostrich is available from https://github.com/usbr/ostrich. The simulation scripts and results of 661 662 the mHM model simulations are publicly available at https://doi.org/10.5281/zenodo.8059198. The source code to quantify the sampling uncertainty in performance metrics (the "gumboot" 663 package) is available at https://github.com/CH-Earth/gumboot. The scripts and results of the 664 665 gumboot-based sampling uncertainty analysis is available at https://doi.org/10.5281/zenodo.8058659 666

Acknowledgements: We acknowledge the financial support for the SPACE project by the
Villum Foundation (http://villumfonden.dk/) through their Young Investigator Program (grant
VKR023443). The first author is supported by NASA program i.e. NNH22ZDA001N-RRNES:
A.24 Rapid Response and Novel Research in Earth Science under the grant number 22RRNES22-0010 and by the Scientific Research Projects Department of Istanbul Technical
University (ITU-BAP) under grant number MDA-2022-43762 and by the National Center for

- High Performance Computing of Turkey (UHeM) under grant number 1007292019.
- 674 **Conflicts of Interest:** "The authors declare no conflict of interest."
- 675 Institutional Review Board Statement: Not applicable.
- 676 **Informed Consent Statement:** Not applicable.
- 677
- 678

### 679 **References**

- Ahmed, K., Sachindra, D. A., Shahid, S., Demirel, M. C., & Chung, E.-S. (2019). Selection of
  multi-model ensemble of general circulation models for the simulation of precipitation
  and maximum and minimum temperature based on spatial assessment metrics. *Hydrology and Earth System Sciences*, 23(11), 4803–4824. https://doi.org/10.5194/hess-23-48032019
- Allen, R. G., Pereira, L. S., Raes, D., & Smith, M. (1998). Crop evapotranspiration -*Guidelines for computing crop water requirements*. FAO Irrigation and drainage paper
  56. Retrieved from http://www.fao.org/docrep/x0490e/x0490e00.htm (accessed at
  16/02/2018)
- Arun, N., Gaw, N., Singh, P., Chang, K., Aggarwal, M., Chen, B., et al. (2021). Assessing the
   Trustworthiness of Saliency Maps for Localizing Abnormalities in Medical Imaging.
   *Radiology: Artificial Intelligence*, 3(6). https://doi.org/10.1148/ryai.2021200267
- Asadzadeh, M., & Tolson, B. (2013). Pareto archived dynamically dimensioned search with
   hypervolume-based selection for multi-objective optimization. *Engineering Optimization*,
   45(12), 1489–1509. https://doi.org/10.1080/0305215X.2012.748046
- Avcuoğlu, M. B., & Demirel, M. C. (2022). Hidrolojik Model Kalibrasyonunda Uydu Tabanlı
  Aylık Buharlaşma ve LAI Verilerinin Kullanılması. *Teknik Dergi*, *33*(6), 13013–13035.
  https://doi.org/10.18400/tekderg.1067466
- Becker, R., Koppa, A., Schulz, S., Usman, M., aus der Beek, T., & Schüth, C. (2019). Spatially
  distributed model calibration of a highly managed hydrological system using remote
  sensing-derived ET data. *Journal of Hydrology*, 577(10), 123944.
  https://doi.org/10.1016/j.jhydrol.2019.123944
- Beven, K. (2023). Benchmarking hydrological models for an uncertain future. *Hydrological Processes*, *37*(5). https://doi.org/10.1002/hyp.14882
- de Boer-Euser, T., Bouaziz, L., De Niel, J., Brauer, C., Dewals, B., Drogue, G., et al. (2017).
  Looking beyond general metrics for model comparison lessons from an international
  model intercomparison study. *Hydrology and Earth System Sciences*, 21(1), 423–440.
  https://doi.org/10.5194/hess-21-423-2017
- Busari, I. O., Demirel, M. C., & Newton, A. (2021). Effect of Using Multi-Year Land Use Land
  Cover and Monthly LAI Inputs on the Calibration of a Distributed Hydrologic Model. *Water*, *13*(11), 1538. https://doi.org/10.3390/w13111538
- Cain, M. K., Zhang, Z., & Yuan, K.-H. (2017). Univariate and multivariate skewness and
   kurtosis for measuring nonnormality: Prevalence, influence and estimation. *Behavior Research Methods*, 49(5), 1716–1735. https://doi.org/10.3758/s13428-016-0814-1
- Clark, M. P., Vogel, R. M., Lamontagne, J. R., Mizukami, N., Knoben, W. J. M., Tang, G., et
  al. (2021). The Abuse of Popular Performance Metrics in Hydrologic Modeling. *Water Resources Research*, 57(9). https://doi.org/10.1029/2020WR029001
- Danapour, M., Fienen, M. N., Højberg, A. L., Jensen, K. H., & Stisen, S. (2021). MultiConstrained Catchment Scale Optimization of Groundwater Abstraction Using Linear
  Programming. *Groundwater*, 59(4), 503–516. https://doi.org/10.1111/gwat.13083
- Dembélé, M., Hrachowitz, M., Savenije, H. H. G., Mariéthoz, G., & Schaefli, B. (2020).
  Improving the Predictive Skill of a Distributed Hydrological Model by Calibration on

- Spatial Patterns With Multiple Satellite Data Sets. *Water Resources Research*, 56(1).
  https://doi.org/10.1029/2019WR026085
- Demirel, M. C., Booij, M. J., & Hoekstra, A. Y. (2013). Effect of different uncertainty sources
   on the skill of 10 day ensemble low flow forecasts for two hydrological models. *Water Resources Research*, 49(7), 4035–4053. https://doi.org/10.1002/wrcr.20294
- Demirel, M. C., Booij, M. J., & Hoekstra, A. Y. (2015). The skill of seasonal ensemble low flow forecasts in the Moselle River for three different hydrological models. *Hydrology and Earth System Sciences*, 19(1), 275–291. https://doi.org/10.5194/hess-19-275-2015
- Demirel, M. C., Mai, J., Mendiguren, G., Koch, J., Samaniego, L., & Stisen, S. (2018).
  Combining satellite data and appropriate objective functions for improved spatial pattern
  performance of a distributed hydrologic model. *Hydrology and Earth System Sciences*,
  22(2), 1299–1315. https://doi.org/10.5194/hess-22-1299-2018
- Dougherty, E., Sherman, E., & Rasmussen, K. L. (2020). Future Changes in the Hydrologic
  Cycle Associated with Flood-Producing Storms in California. *Journal of Hydrometeorology*, 21(11), 2607–2621. https://doi.org/10.1175/JHM-D-20-0067.1
- Efford, N. (2000). Digital Image Processing: A Practical Introduction Using JavaTM. Pearson
   Education. Slate.
- Efron, B. (1992). Jackknife-After-Bootstrap Standard Errors and Influence Functions. *Journal*of the Royal Statistical Society: Series B (Methodological), 54(1), 83–111.
  https://doi.org/10.1111/j.2517-6161.1992.tb01866.x
- Fréchet, M. M. (1906). Sur quelques points du calcul fonctionnel. *Rendiconti Del Circolo Matematico Di Palermo*, 22(1), 1–72. https://doi.org/10.1007/BF03018603
- Freedman, D., & Diaconis, P. (1981). On the histogram as a density estimator:L 2 theory. *Zeitschrift Für Wahrscheinlichkeitstheorie Und Verwandte Gebiete*, 57(4), 453–476.
  https://doi.org/10.1007/BF01025868
- Gaur, S., Singh, B., Bandyopadhyay, A., Stisen, S., & Singh, R. (2022). Spatial pattern-based
   performance evaluation and uncertainty analysis of a distributed hydrological model.
   *Hydrological Processes*, *36*(5). https://doi.org/10.1002/hyp.14586
- Gupta, H. V., Kling, H., Yilmaz, K. K., & Martinez, G. F. (2009). Decomposition of the mean
  squared error and NSE performance criteria: Implications for improving hydrological
  modelling. *Journal of Hydrology*, 377(1–2), 80–91.
  https://doi.org/10.1016/j.jhydrol.2009.08.003
- Hargreaves, G. H., & Samani, Z. A. (1985). Reference Crop Evapotranspiration from
  Temperature. *Applied Engineering in Agriculture*, 1(2), 96–99.
  https://doi.org/10.13031/2013.26773
- Hausdorff, F. (1914). Bemerkung über den Inhalt von Punktmengen. *Mathematische Annalen*,
  758 75(3), 428–433.
- Hossain, M. K., & Meng, Q. (2020). A thematic mapping method to assess and analyze
  potential urban hazards and risks caused by flooding. *Computers, Environment and Urban Systems*, 79, 101417. https://doi.org/10.1016/j.compenvurbsys.2019.101417
- Immerzeel, W. W., & Droogers, P. (2008). Calibration of a distributed hydrological model
  based on satellite evapotranspiration. *Journal of Hydrology*, 349(3–4), 411–424.

- 764 https://doi.org/10.1016/j.jhydrol.2007.11.017
- Knoben, W. J. M., & Spieler, D. (2022). Teaching hydrological modelling: illustrating model
  structure uncertainty with a ready-to-use computational exercise. *Hydrology and Earth System Sciences*, 26(12), 3299–3314. https://doi.org/10.5194/hess-26-3299-2022
- Knoben, W. J. M., Freer, J. E., & Woods, R. A. (2019). Technical note: Inherent benchmark or
   not? Comparing Nash-Sutcliffe and Kling-Gupta efficiency scores. *Hydrology and Earth System Sciences Discussions*, (July), 1–7. https://doi.org/10.5194/hess-2019-327
- Kumar, R., Samaniego, L., & Attinger, S. (2013). Implications of distributed hydrologic model
   parameterization on water fluxes at multiple scales and locations. *Water Resources Research*, 49(1), 360–379. https://doi.org/10.1029/2012WR012195
- López, P., Sutanudjaja, E. H., Schellekens, J., Sterk, G., & Bierkens, M. F. P. (2017).
  Calibration of a large-scale hydrological model using satellite-based soil moisture and
  evapotranspiration products. *Hydrology and Earth System Sciences*, 21(6), 3125–3144.
  https://doi.org/10.5194/hess-21-3125-2017
- Martinez-Villalobos, C., Neelin, J. D., & Pendergrass, A. G. (2022). Metrics for Evaluating
   CMIP6 Representation of Daily Precipitation Probability Distributions. *Journal of Climate*, *35*(17), 5719–5743. https://doi.org/10.1175/JCLI-D-21-0617.1
- Matott, L. Shawn. (2004). OSTRICH: an Optimization Software Tool, Documentation and
   User's Guide, Version 17.12.19. Retrieved from https://github.com/usbr/ostrich
- Matott, L.S. (2017). OSTRICH: an Optimization Software Tool, Documentation and User's
   Guide. University at Buffalo Center for Computational Research, Version 17, 79.
- Monteith, J. L. (1965). EVAPORATION AND ENVIRONMENT. Symposia of the Society for
   *Experimental Biology*, 19, 205–234.
- Nash, J. E., & Sutcliffe, J. V. (1970). River flow forecasting through conceptual models part I
  A discussion of principles. *Journal of Hydrology*, 10(3), 282–290.
  https://doi.org/10.1016/0022-1694(70)90255-6
- Nilsson, J., & Akenine-Möller, T. (2020). Understanding SSIM. Retrieved from https://arxiv.org/abs/2006.13846
- Norman, J. M., Kustas, W. P., & Humes, K. S. (1995). Source approach for estimating soil and
   vegetation energy fluxes in observations of directional radiometric surface temperature.
   *Agricultural and Forest Meteorology*, 77(3–4), 263–293. https://doi.org/10.1016/0168 1923(95)02265-Y
- Odusanya, A. E., Schulz, K., & Mehdi-Schulz, B. (2022). Using a regionalisation approach to
  evaluate streamflow simulated by an ecohydrological model calibrated with global land
  surface evaporation from remote sensing. *Journal of Hydrology: Regional Studies*, 40,
  101042. https://doi.org/10.1016/j.ejrh.2022.101042
- 800 Onyutha, C. (2022). A hydrological model skill score and revised R-squared. *Hydrology* 801 *Research*, 53(1), 51–64. https://doi.org/10.2166/nh.2021.071
- Pearson, K. (1905). The problem of the random walk. *Nature*, 72(1865), 294.
- Priestley, C. H. B., & Taylor, R. J. (1972). On the Assessment of Surface Heat Flux and
  Evaporation Using Large-Scale Parameters. *Monthly Weather Review*, 100(2), 81–92.

- 805 https://doi.org/10.1175/1520-0493(1972)100<0081:OTAOSH>2.3.CO;2
- Rakovec, O., Kumar, R., Mai, J., Cuntz, M., Thober, S., Zink, M., et al. (2016). Multiscale and
   Multivariate Evaluation of Water Fluxes and States over European River Basins. *Journal of Hydrometeorology*, *17*(1), 287–307. https://doi.org/10.1175/JHM-D-15-0054.1
- Rientjes, T. H. M., Muthuwatta, L. P., Bos, M. G., Booij, M. J., & Bhatti, H. A. (2013). Multi-variable calibration of a semi-distributed hydrological model using streamflow data and satellite-based evapotranspiration. *Journal of Hydrology*, 505, 276–290. https://doi.org/10.1016/j.jhydrol.2013.10.006
- Rubner, Y., Tomasi, C., & Guibas, L. J. (1998). A metric for distributions with applications to
  image databases. In *Sixth International Conference on Computer Vision (IEEE Cat. No.98CH36271)* (pp. 59–66). Narosa Publishing House.
  https://doi.org/10.1109/ICCV.1998.710701
- Samaniego, L., Kumar, R., & Attinger, S. (2010). Multiscale parameter regionalization of a
  grid-based hydrologic model at the mesoscale. *Water Resources Research*, 46(5),
  W05523. https://doi.org/10.1029/2008WR007327
- Samaniego, L., Brenner, J., Craven, J., Cuntz, M., Dalmasso, G., Demirel, C. M., et al. (2021,
  July 21). The mesoscale Hydrologic Model mHM v5.11.2.
  https://doi.org/10.5281/ZENODO.5119952
- Schneider, R., Henriksen, H. J., & Stisen, S. (2022). A robust objective function for calibration
  of groundwater models in light of deficiencies of model structure and observations. *Journal of Hydrology*, *613*, 128339. https://doi.org/10.1016/j.jhydrol.2022.128339
- Scott, D. W. (1979). On optimal and data-based histograms. *Biometrika*, 66(3), 605–610.
  https://doi.org/10.1093/biomet/66.3.605
- Sirisena, T. A. J. G., Maskey, S., & Ranasinghe, R. (2020). Hydrological Model Calibration
  with Streamflow and Remote Sensing Based Evapotranspiration Data in a Data Poor
  Basin. *Remote Sensing*, *12*(22), 3768. https://doi.org/10.3390/rs12223768
- Stisen, S., Koch, J., Sonnenborg, T. O., Refsgaard, J. C., Bircher, S., Ringgaard, R., & Jensen,
  K. H. (2018). Moving beyond run-off calibration-Multivariable optimization of a surfacesubsurface-atmosphere model. *Hydrological Processes*, 32(17), 2654–2668.
  https://doi.org/10.1002/hyp.13177
- Sturges, H. A. (1926). The Choice of a Class Interval. *Journal of the American Statistical Association*, 21(153), 65–66. https://doi.org/10.1080/01621459.1926.10502161
- Swain, M. J., & Ballard, D. H. (1991). Color indexing. *International Journal of Computer Vision*, 7(1), 11–32. https://doi.org/10.1007/BF00130487
- Thober, S., Cuntz, M., Kelbling, M., Kumar, R., Mai, J., & Samaniego, L. (2019). The multiscale routing model mRM v1.0: simple river routing at resolutions from 1 to 50 km. *Geoscientific Model Development*, *12*(6), 2501–2521. https://doi.org/10.5194/gmd-12-2501-2019
- Thoya, P., Maina, J., Möllmann, C., & Schiele, K. S. (2021). AIS and VMS Ensemble Can
  Address Data Gaps on Fisheries for Marine Spatial Planning. *Sustainability*, *13*(7), 3769.
  https://doi.org/10.3390/su13073769
- 846 Tolson, B. A., & Shoemaker, C. A. (2007). Dynamically dimensioned search algorithm for

- computationally efficient watershed model calibration. *Water Resources Research*, 43(1).
  https://doi.org/10.1029/2005WR004723
- Wakigari, S. A., & Leconte, R. (2023). Assessing the Potential of Combined SMAP and InSitu Soil Moisture for Improving Streamflow Forecast. *Hydrology*, *10*(2), 31.
  https://doi.org/10.3390/hydrology10020031
- Wang, Z., Bovik, A. C., Sheikh, H. R., & Simoncelli, E. P. (2004). Image Quality Assessment:
  From Error Visibility to Structural Similarity. *IEEE Transactions on Image Processing*, *13*(4), 600–612. https://doi.org/10.1109/TIP.2003.819861
- Wiederholt, R., Paudel, R., Khare, Y., Davis, S. E., Melodie Naja, G., Romañach, S., et al.
  (2019). A multi-indicator spatial similarity approach for evaluating ecological restoration
  scenarios. *Landscape Ecology*, *34*(11), 2557–2574. https://doi.org/10.1007/s10980-01900904-w
- Yoo, S. B. M., Tu, J. C., Piantadosi, S. T., & Hayden, B. Y. (2020). The neural basis of
   predictive pursuit. *Nature Neuroscience*, 23(2), 252–259. https://doi.org/10.1038/s41593 019-0561-6
- Zink, M., Mai, J., Cuntz, M., & Samaniego, L. (2018). Conditioning a Hydrologic Model Using
   Patterns of Remotely Sensed Land Surface Temperature. *Water Resources Research*,
   54(4), 2976–2998. https://doi.org/10.1002/2017WR021346

# Benchmarking multi-component spatial metrics for hydrologic model calibration using MODIS AET and LAI products

# 3 E. B. Yorulmaz<sup>1</sup>, E. Kartal<sup>1</sup> and M. C. Demirel<sup>1</sup>

<sup>1</sup> Department of Civil Engineering, Istanbul Technical University, 34467 Maslak, Istanbul,
 Turkey

6

Corresponding author: Eymen Berkay Yorulmaz (<u>yorulmaz21@itu.edu.tr</u>, ORCID 0000-0003-3370-9465)

# 9 Key Points:

- Newly proposed spatial metrics offer significant improvements in discriminating
   between two raster maps
- Selecting appropriate spatial metric proved to be very crucial even for the global search algorithms
- Sampling uncertainty in metrics increases with newly added components
- 15

# 16 Abstract

17 SPAtial EFficiency (SPAEF) metric is one of the most thoroughly metrics in hydrologic community. In this study, our aim is to improve SPAEF by replacing the histogram match 18 19 component with other statistical indices, i.e. kurtosis and earth mover's distance, or by adding 20 a fourth or fifth component such as kurtosis and skewness. The existing spatial metrics i.e. 21 SPAtial Efficiency (SPAEF), Structural Similarity (SSIM) and Spatial Pattern Efficiency 22 Metric (SPEM) were compared with newly proposed metrics to assess their converging performance. The mesoscale Hydrologic Model (mHM) of the Moselle River is used to 23 24 simulate streamflow (Q) and actual evapotranspiration (AET). The two-source energy balance 25 (TSEB) AET during the growing season is used as monthly reference maps to calculate the 26 spatial performance of the model. The Moderate Resolution Imaging Spectroradiometer (MODIS) based Leaf area index (LAI) is utilized by the mHM via pedo-transfer functions and 27 28 multi-scale parameter regionalization approach to scale the potential ET. In addition to the real 29 monthly AET maps, we also tested these metrics using a synthetic true AET map simulated 30 with a known parameter set for a randomly selected day. The results demonstrate that the newly 31 developed four-component metric i.e. SPAtial Hybrid 4 (SPAH4) slightly outperform 32 conventional three-component metric i.e. SPAEF (3% better). However, SPAH4 significantly 33 outperforms the other existing metrics i.e. 40% better than SSIM and 50% better than SPEM. 34 We believe that other fields such as remote sensing, change detection, function space optimization and image processing can also benefit from SPAH4. 35

- 36 Keywords: mHM, model calibration, spatial pattern, SPAEF, MODIS, TSEB
- 37

### 38 1 Introduction

39 Distributed hydrologic models have a crucial role in creating digital twin of the water cycle in nature by revealing physical mechanisms and process interactions. After identifying the best 40 41 parameter set through calibration, these models are used to conduct robust numerical experiments assessing climate change impacts (Beven, 2023) or land use land cover change 42 43 impacts on model output fluxes such as runoff (Busari et al., 2021), groundwater recharge, soil moisture and actual evapotranspiration (AET). A skillful model enables decision-makers to 44 plan for and respond to water-related extremes such as hydrological droughts and floods. 45 46 Accuracy of the model results depends on the success of identifying best combination of the 47 parameters since calibration process helps us reduce discrepancies in model physics. Demirel et al. (Demirel et al., 2018) showed that using only streamflow hydrograph performance as 48 49 objective function diminishes the AET patterns simulated by the model. However, incorporating satellite based remotely sensed AET into the multi-objective calibration 50 51 framework that has already streamflow, surprisingly improves both water balance and AET 52 performance of the model. Other studies benefitted from land surface temperature (Zink et al., 2018), soil moisture (López et al., 2017; Wakigari & Leconte, 2023), AET (Avcuoğlu & 53 Demirel, 2022; Gaur et al., 2022; Odusanya et al., 2022; Sirisena et al., 2020) and groundwater 54 55 (Danapour et al., 2021; Stisen et al., 2018) in hydrologic model calibration.

56 In other words, hydrologic model calibration is essential for ensuring the validity and reliability of model predictions i.e. of most important for water management and decision-making 57 processes. However, the robustness of hydrologic model calibration heavily relies on how the 58 59 model is guided in the solution space via the performance metrics (de Boer-Euser et al., 2017; Knoben et al., 2019; Martinez-Villalobos et al., 2022; Onyutha, 2022; Schneider et al., 2022). 60 61 If the metric is too loose (tolerant) or prone to the sampling uncertainty (Clark et al., 2021), the 62 calibration process can stop quickly in the local minima while the modeler searches for the best global solution. The key point of the modelling chain is the selection of appropriate metric. 63 64 Our study focuses on development of a novel metric with least tolerance (highest 65 discrimination skill) based on benchmarking existing metrics in evaluating the similarity of two raster maps. We are particularly interested in multi-component bias-insensitive spatial 66 67 metrics for pattern comparison. Thus, bias sensitive temporal metrics used for water balance are not within the scope of this study. 68

69 The use of multi-component spatial metrics in hydrologic model calibration is an important advancement in the field of water resource management and resource allocation. The multi-70 71 component metrics provides a more nuanced evaluation of model performance compared to 72 traditional single-component metrics e.g. mean absolute error and coefficient of determination. 73 The adoption of these metrics allows for a more comprehensive understanding of the 74 hydrologic system and its spatial variability, which is critical for informed decision-making. 75 These metrics differ from single-component metrics in that they consider multiple components of the hydrological system, rather than just one component. By providing a more 76 77 comprehensive evaluation of the hydrologic system, multi-component metrics help to identify 78 areas where models can be improved. For spatial metrics, the added level of complexity 79 provided by multi-component metrics offers a more robust evaluation of model performance, 80 providing a better understanding of the spatial variability of the hydrologic system.

In recent years, remote sensing data from satellites, such as Moderate Resolution Imaging Spectroradiometer (MODIS) products, have become commonly used in hydrologic model calibration since this product provides estimates of AET from vegetation, which is a key component and major water loss in the hydrologic cycle (Becker et al., 2019; Rientjes et al., 2013). On one hand, it serves to better represent the cell-to-cell hydrological dynamics and 86 diversity in the basin also allows for a more detailed understanding of the water budget at the 87 land surface and helps to better quantify the water requirements of vegetation. On the other hand, the MODIS Leaf Area Index (LAI), product provides information about the leaf area 88 index, which is a measure of the amount of vegetation cover in an area. This information is 89 90 essential for understanding how vegetation influences the water cycle by affecting factors such as precipitation, evapotranspiration, and runoff. In this study, we use LAI to dynamically scale 91 92 the PET input to the model to improve AET performance and present a comprehensive 93 benchmarking of multi-component spatial metrics using MODIS-LAI and TSEB AET 94 products, to assess their potential for calibration (Immerzeel & Droogers, 2008).

95 There are various performance metrics in hydrology. The Nash-Sutcliffe Efficiency (NSE) and 96 Kling-Gupta Efficiency (KGE) are the most widely recognized performance metrics used in 97 evaluating and calibrating rainfall-runoff models. These two metrics have been instrumental in 98 advancing our understanding of hydrological processes and improving the performance of 99 hydrologic models (Gupta et al., 2009; Nash & Sutcliffe, 1970). They have paved the way for 100 the development of more advanced and sophisticated performance evaluation techniques. Despite the sampling uncertainty inherited in these metrics (Clark et al., 2021), NSE and KGE 101 102 continue to be widely accepted in the hydrology community due to their simplicity and 103 effectiveness in evaluating model performance. Many of the newer metrics that have been 104 introduced in recent years have been inspired by and built upon the foundation established by 105 NSE and KGE. The conventional model calibration relies on using flow-oriented temporal 106 metrics, such as the NSE and KGE. However, these metrics have a limitation as they lack 107 spatial considerations and are prone to the sampling uncertainty. This has driven the need for 108 development of intolerant spatial performance metrics which can better evaluate and improve 109 the spatial accuracy of a hydrologic model. Spatial-pattern-oriented SPAtial Efficiency (SPAEF) metric developed by Demirel et al. (Demirel et al., 2018) builds upon the strength of 110 111 KGE and incorporates new idea of distribution comparison via histogram overlap index. It is designed as a multi-component metric specifically suited for comparing spatial patterns of two 112 raster maps, with its three main data properties being co-location, variation, and distribution. 113 Although SPAEF was primarily developed for hydrologic community, it has been used in many 114 115 different disciplines such as atmospheric circulation modeling (Ahmed et al., 2019), flood risk 116 analysis (Hossain & Meng, 2020), function space optimization, fisheries (Thoya et al., 2021) and neuroscience (Yoo et al., 2020). In these studies, SPAEF has been tested and proven to be 117 118 robust and easy to interpret due to its three distinct and complementary components of 119 correlation, variance and histogram matching. Following the multi-component structure idea, we present new metrics in this study to improve SPAEF by adding fourth of fifth new 120 121 components or replacing histogram match with other components. Using this approach, we 122 aimed for reducing uncertainty in the new metric and make it sharp (discriminant) when 123 evaluating patterns on two raster maps whether they are similar or not.

124 In recent literature, there has been attempts to revise SPAEF component i.e. Spatial Pattern 125 Efficiency Metric (SPEM) (Dembélé et al., 2020). Similar to SPAEF, it has been proposed as 126 a bias-insensitive and multi-component spatial pattern-oriented metric using satellite remote 127 sensing data. Structural Similarity index (SSIM) is another pattern-oriented metric, it stands out with its spatial structure (Nilsson & Akenine-Möller, 2020; Wang et al., 2004). It was 128 129 proposed by Wang et al. (Wang et al., 2004) for image quality assessment and has been used 130 in different studies such as medical imaging, ecological restoration, and change detection in 131 the hydrological cycles and remote sensing images (Arun et al., 2021; Dougherty et al., 2020; Wiederholt et al., 2019). Knoben et al. (Knoben et al., 2019) compared NSE and KGE metrics 132 133 and argued that instead of relying directly on the KGE value, the components should be 134 analyzed in depth, even the weighting of the components. A study analyzing sampling 135 uncertainty in popular performance metrics in hydrologic modeling highlighted that the KGE can be heavily influenced by just a few data points (Clark et al., 2021). A study on the 136 hydrological model skill score compared metrics with different forms of correlation and 137 138 measures of variability, claiming the term covariance is more appropriate for evaluation 139 (Onyutha, 2022). Another recent study, based on the largest residuals, focused on reducing the largest errors, and argued that metrics should be less sensitive to errors and more sensitive to 140 141 bias (Schneider et al., 2022). The publication (Martinez-Villalobos et al., 2022) compared 142 metrics for evaluating precipitation probability distributions by comparing climate model 143 simulation data with real platform satellite data, therefore they showed the importance of 144 probability distribution functions. A study from the Netherlands (de Boer-Euser et al., 2017) 145 stated that strong components can be included in different metrics rather than considering a 146 single general metric for model comparison.

147 The existing spatial metrics aimed for the best convergence using terms such as correlation, 148 variation, histogram intersection, and root mean square error. However, kurtosis has hitherto 149 been an underrated term for spatial performance, and a four-component spatial-pattern-oriented 150 metric also does not exist for the hydrologic model calibration. We used the kurtosis ratio by 151 including it as a new component for the first time in this study in order to achieve the best 152 spatial convergence and fit. With the addition of a new component, the weighting by which the 153 components affect the value has also changed. By revealing the effect of kurtosis on spatial 154 performance, we developed a new four-component metric that does not require user input.

155 We aim to investigate the best potential to use multi-component spatial metrics in hydrological 156 model calibration, by proposing a new multi-component spatial metric that especially includes 157 the kurtosis component and benchmarking it to existing multi-component spatial metrics. The 158 primary purpose of this study is to evaluate the performance of the hydrological model using 159 multicomponent spatial metrics and to determine the potential impact on model accuracy and 160 precision. In addition, this study aims to identify the most effective combination of spatial metrics for hydrological model calibration and to develop a framework for future work in this 161 area. A large number of metrics in the literature creates confusion and difficulty for users to 162 163 choose from, so we compared metrics in this study to look for the most successful one to put a 164 stop to metric redundancy. Addressing these goals, this study aims to contribute to ongoing research efforts to improve the accuracy and reliability of hydrological models. 165

The accuracy of the analysis has been increased by comparing model predictions with real 166 platforms. It is aimed to improve the convergence between observed and simulated maps by 167 168 using two-source energy balance (TSEB) model's AET data. The MODIS-LAI data were used 169 both to correct the PET and to represent the vegetation dynamics of the Moselle basin. We utilize a spatially distributed mesoscale Hydrologic Model (mHM) with it features pedo-170 171 transfer functions for LAI data and a Multiscale Parameter Regionalization (MPR) approach to scale the potential ET (Kumar et al., 2013; Samaniego et al., 2010). We tested our framework 172 173 in three different cases to provide comprehensive outlook to the calibrations i.e. 100 iterations 174 were applied in the first case and 1000 iterations in the second case, so the effect of the number 175 of iterations was also assessed. In the third case, reproducibility was achieved by analyzing the randomly selected synthetic map. OSTRICH software (L. Shawn Matott, 2004; L.S. Matott, 176 177 2017) was used as the calibration tool and Parallel Dynamically Dimensioned Search 178 Algorithm (PDDS) was used as the calibration algorithm (Asadzadeh & Tolson, 2013). The 179 combined SPAEF value of the growing season was used as the main objective function for ET, 180 and the KGE was presented for discharge (Q) in addition. We developed multiple metrics with 181 different components and different component numbers, trying to increase the effectiveness (sharpness) of each component on convergence performance. We made an elaborated 182

183 comparison between the existing performance metrics in the literature and the newly developed 184 metrics based on SPAEF. As a result of the rigorous assessment of metrics, we identified not 185 only the superior but also new metric. The strongest aspect of this new metric is the added 186 kurtosis component.

### 187 2 Study area and data

### 188 2.1 Study area

189 The study area is the Moselle River basin, the largest part of the Rhine River basin, of which it is one of the main tributaries, characterized by diverse landforms (Figure 1). The origin of the 190 191 river from the Vosges Mountains before the interterritorial transfer from France to enter 192 Germany and Luxembourg. Furthermore, at the triangle where Germany, France and 193 Luxembourg meet, the Moselle River becomes the borderline between Germany and Luxembourg for 36 km. Also, it has a surface area of approximately 27262 km<sup>2</sup> and a length 194 195 of 545 km. Whereas, land use in the basin includes forestry, agriculture and cattle breeding in 196 the mountains and hillslopes, winegrowing on vineyards of sunny valley slopes. Moreover, the 197 altitude varies from 59 to 1326 m, with an average altitude of around 340 m (Demirel et al., 198 2013). In addition to having 26 sub-basins with surface areas varying from 102 to 3353 km<sup>2</sup>, 199 the river flow is organized by different dams, dikes, powerplants and locks such as the Trier 200 Dam, Koblenz Dam and Detzem Lock. The outlet discharge at Cochem station, located between Trier and Koblenz, varies from 14 m<sup>3</sup>/s in dry summers to a maximum of 4000 m<sup>3</sup>/s 201 202 during winter floods, with a mean discharge of around  $315 \text{ m}^3/\text{s}$  (Demirel et al., 2015).



203 204

Figure 1. DEM, land cover and AET characteristics of Mosel River basin.

An average pattern of satellite-based actual evapotranspiration for July (average of all years from 2002 to 2014) is presented to illustrate the interaction between DEM and land cover characteristics that generate the land surface flux patterns.

208 2.2 Satellite data

MODIS has a vital role in obtaining the satellite-based data used in this study, is an essential sensor aboard the Terra (EOS AM) and Aqua (EOS PM) satellites for the earth and climate measurements at a spatial resolution of approximately  $1 \text{ km} \times 1 \text{ km}$ . It provides terrestrial, atmospheric and thalassic data and a view of the entire Earth's surface for large and diverse user communities around the world. In this study, TSEB based AET is used as reference spatial

- 214 patterns (Allen et al., 1998; Norman et al., 1995). TSEB is an energy balance model using the
- energy flux principle by separating into two-layer, vegetation and soil.

216 The water limited growing season was chosen as the analysis period because it avoids climate

gradient on the AET patterns emphasizing vegetation dynamics instead of wet soil conditionsi.e. AET that is equal to the PET. All remote-sensing-based AET data were converted to long

term monthly mean data during the growing season across all years for the model calibration period (2002–2014). In what follows, three-monthly mean periods were obtained with a total

- of three-term between March and November, i.e. March-April-May (MAM), June-July-August
- 222 (JJA), and September-October-November (SON), representing AET under cloud-free
- conditions. We will attribute these AET maps as reference observations, although they are estimates from an energy balance model based on satellite observations and not pure
- 225 observations.
- Table 1. Overview of morphological and meteorological data used as input for mHM (Rakovecet al., 2016).

Variable	Description	Spatial resolution (degrees)	Source
Q (daily)	Streamflow	Point	GRDC
P (daily)	Precipitation	0.0625	E-OBS
PET (daily)	Potential evapotranspiration based on Hargreaves and Samani (Hargreaves & Samani, 1985)	0.0625	E-OBS
Tavg	Average air temperature	0.0625	E-OBS
LAI	Fully distributed 12-monthly values based on 8- day time-varying leaf area index (LAI) dataset	0.001953125	MODIS
Land cover	Forest, agriculture and urban	0.001953125	MODIS
DEM-related data	Slope, aspect, flow accumulation and direction	0.001953125	SRTM
Geology class	Two main geological formations	0.001953125	ESD UFZ – Leipzig (Rakovec et al., 2016)
Soil class	Fully distributed soil texture data	0.001953125	HWSD

228 GRDC - Global Runoff Data Centre, E-OBS - The gridded observational dataset from Copernicus, MODIS - Moderate

229 Resolution Imaging Spectroradiometer, SRTM – Shuttle Radar Topography Mission, ESD – European Soil Database, HWSD –

230 Harmonized World Soil Database

### 231 **3 Hydrological model**

232 This research utilizes the mesoscale Hydrologic Model (mHM) v.5.11.2 (Samaniego et al., 233 2021) which is a grid-based spatially distributed model it features pedo-transfer functions and 234 MPR (Kumar et al., 2013; Samaniego et al., 2010; Thober et al., 2019). Another feature of 235 mHM is the use of leaf area index (LAI) data not only for calculating interception loss but also 236 for dynamically scaling PET (Demirel et al., 2018). With these unique features, it is more 237 flexible than other existing hydrologic models in line with the purpose of this study. The model 238 features 69 adjustable global parameters that can be optimized during the calibration process 239 (Demirel et al., 2018). The model works on the basis of water balance rather than energy 240 balance and provides various physically meaningful spatial outputs, fluxes and states as 241 simulating major elements of the hydrologic processes, i.e. soil moisture dynamics, 242 interception, infiltration, evapotranspiration, snow accumulation and melting, groundwater storage, seepage, surface runoff and others. 243

244 The basic data for the running mHM can be classified into meteorological data, morphological 245 data, land cover data and gauge streamflow data. Table 1 shows a summary of the data used in 246 mHM setup provided by Rakovec et al. (Rakovec et al., 2016). As seen in the table, mHM can 247 handle different spatial resolutions of meteorological data and morphological data since it has internal upscaling and downscaling subroutines. At this point, the Multi-Scale Parameter 248 249 Regionalization technique comes into play and enables user to map calibrated parameters to 250 the simulated grids with pedo-transfer functions. This approach prevents uniform parameter 251 fields and protects sub-grid heterogeneity of the fluxes. In other models, every parameter gets the same value in the entire sub-basin or in each hydrologic response units resulting in uniform 252 253 flux results for the same domain.

254 The meteorological model inputs are precipitation, average air temperature and potential 255 evapotranspiration (PET). In our study, PET was direct input to the mHM and estimated outside 256 with Hargreaves-Samani (Hargreaves & Samani, 1985) method using additional temperature 257 data. All meteorological data are obtained from E-OBS at daily resolution, originally at 10-20 258 km. The morphological variables are digital elevation model (DEM), soil maps with textural 259 features, geological maps including specific yield, permeability and aquifer thickness. In 260 addition to characterizing the morphology of the basin, DEM masks the grid cells with the 261 basin boundaries to eliminate no-data parts. All morphological data are prepared at 262 0.001953125 degrees ( $\sim$ 200 m  $\times$  200 m) scale. The model hydrology is evaluated at 0.015625 degrees (~2x2 km) spatial resolution and daily time step. Lastly, monthly leaf area index (LAI) 263 264 maps are used to represent the vegetation dynamics for both interception calculation and PET 265 correction for the entire period (2002–2014). Four years of model warm-up period (1998–2001) 266 is used. Observed daily streamflow (Q) data at Cochem (station #6336050), provided by the Global Runoff Data Centre (GRDC), Koblenz (Germany), is used to calibrate water balance in 267 the basin. 268

### 269 4 Methods

270 In this study, we tested nine different spatial metrics i.e. two of them are existing metrics, and 271 seven of them are newly developed based on SPAEF (Table 2). To evaluate the effect of 272 number of iterations, calibrations were pursued with either 100 or 1000 maximum iterations. 273 Besides, synthetically created AET maps using mHM and a pre-defined parameter set are 274 utilized to mimic a "hide and seek" case. This is crucial to test the guidance performance of the 275 metrics in the multi-dimensional solution space to find the hided (perfect) solution within 1000 276 iterations since search algorithms, i.e. ParaPADDS algorithm herein, require a metric to 277 evaluate model results at every iteration.

### 278 4.1 Objective Functions

Multi-component structure of our metrics was inspired by the Kling–Gupta efficiency (Gupta et al., 2009). KGE is one of the most used metrics in the hydrologic modelling to evaluate streamflow performance. As shown in Eq.(1), it has three components, i.e., correlation, variability and bias.

$$KGE = 1 - \sqrt{\left(\alpha_Q - 1\right)^2 + \left(\beta_Q - 1\right)^2 + \left(\gamma_Q - 1\right)^2}$$

$$\alpha_Q = \rho(o, s), \beta_Q = \frac{\sigma_S}{\sigma_0} \text{ and } \gamma_Q = \frac{\mu_S}{\mu_0}$$
(1)

where  $\alpha_Q$  is the Pearson correlation coefficient between the observed (o) and the simulated (s) discharge time series,  $\beta_Q$  is the relative variability based on the ratio of standard deviation in simulated and observed values and  $\gamma_Q$  is the bias fraction which is normalized by the standard deviation of the observed data.

Table 2 shows the summary of SPAEF based metrics. For brevity, we used Eq. (2) as formula template i.e. a generic formulation type that encompasses in the number and content of components. The excessed style in Eq (2) includes all metrics form with various components.

METRIC = 
$$1 - \sqrt{(\alpha - 1)^2 + (\beta - 1)^2 + (\gamma - 1)^2 + (\kappa - 1)^2 + (\delta - 1)^2}$$
 (2)

Matuia	Components								
	α	β	γ	к	δ				
SPAtial Efficiency (SPAEF)	$\rho(o,s)$	$\frac{\sigma_o}{\mu_o}/\frac{\sigma_s}{\mu_s}$	$\frac{\sum_{j=1}^{n} \min(K_j, L_j)}{\sum_{j=1}^{n} K_j}, n=100 \text{ fixed}$	none	none	Eq. (3)			
SPAtial EFficiency Prime (SPAEF')	same as SPAEF	same as SPAEF	same as SPAEF except for dynamic <i>n</i> i.e. number of bins $n = floor\{\sqrt{length(o)}\}$	none	none	Eq. (4)			
SPAtial Count Density Efficiency (SPACD)	same as SPAEF	same as SPAEF	$\frac{\sum_{j=1}^{n} \min(K_j, L_j)}{\sum_{j=1}^{n} K_j} \left( v_n = c_n / w_n \right)$	none	none	Eq. (5)			
SPAtial Hybrid 4 Efficiency (SPAH4)	same as SPAEF	same as SPAEF	same as SPAEF'	Kurt(s) Kurt(o)	none	Eq. (6)			
SPAtial Kurtosis Efficiency (SPAK)	same as SPAEF	same as SPAEF	none	same as SPAH4	none	Eq. (7)			
SPAtial Hybrid 5 Efficiency (SPAH5)	same as SPAEF	same as SPAEF	same as SPAEF'	same as SPAH4	Skew(s) Skew(o)	Eq. (8)			
SPAtialHistogramEqualizationEfficiency(SPAHE)Efficiency	same as SPAEF	same as SPAEF	$\frac{\sum_{j=1}^{n} min(K_j, L_j)}{\sum_{j=1}^{n} K_j}$	none	none	Eq. (9)			
SPAtial Movers' Distance Efficiency ( <b>SPAMD</b> )	same as SPAEF	same as SPAEF	$\frac{\sum_{i=1}^{K} \sum_{i=1}^{L} f_{i,j} d_{i,j}}{\sum_{i=1}^{K} \sum_{i=1}^{L} f_{i,j}}$	none	none	Eq. (10)			
Spatial Pattern Efficiency Metric ( <b>SPEM</b> )	$1 - \frac{6\sum_1^n d^2}{n(n^2 - 1)}$	same as SPAEF	$1 - E_{RMS}(Z_{X_s}, Z_{X_o})$	none	none	Eq. (11)			

 Table 2. SPAEF based metrics used as objective functions.

291

293 SPAEF is the seed of our newly proposed metrics as our aim is to sharpen SPAEF. In other 294 words, we intend to improve its discriminating power while judging whether two maps are 295 similar or not. SPAEF uses a multi-component structure of the KGE metric. In Eq. (3),  $\alpha$  is the 296 Pearson correlation coefficient between the observed (o) and simulated (s) pattern,  $\beta$  is the 297 fraction of the coefficient of variation representing spatial variability and  $\gamma$  is the histogram 298 intersection, which based on z-scores, for the given histogram K of the observed pattern and 299 the histogram L of the simulated pattern, each containing n bins (Swain & Ballard, 1991). The 300 SPAEF can have a value between  $-\infty$  and 1, where a value closer to 1 indicates highest spatial 301 similarity between the observations and model simulations (Koch et al., 2018).

As a result of various adjustments and improvements made in the SPAEF components, new
 metrics were proposed and tested i.e. SPAEF', SPACD, SPAH4, SPAK, SPAH5, SPAMD, and
 SPAHE. We included two popular metrics, SPEM and SSIM into benchmark.

First improvement in SPAEF is changing user defined the number of bins to an automated n based on the number of elements (grids) in the raster map (see Eq (3)). We introduced a simple approach i.e. the square root of the length of the observed data as n = $floor{\sqrt{length(o)}}$  although there are different methods for the same purpose (Freedman & Diaconis, 1981; Scott, 1979; Sturges, 1926). This slightly new version of the SPAEF is presented as SPAEF-Prime (SPAEF') as shown in Eq (4). Unlike the standard version, the SPAEF' does not require any user-defined inputs now.

Eq (5) shows Spatial Count Density Efficiency (SPACD) which has a different type of normalization based on count density approach in the calculation of the histogram intersection component. While the first two components remain constant as in SPAEF' the calculation of n in the gamma component has changed. This approach uses count or frequency scaled by the width of the bin  $v_n = c_n / w_n$ ,  $v_n$  is the bin value,  $c_n$  is the number of elements in the bin and  $w_n$  is the width of the bin, respectively.

Eq (6) shows SPAtial Hybrid 4 Efficiency (SPAH4) which is a four-component metric obtained 318 319 by adding kurtosis i.e. a fundamental statistical property of distributions to the SPAEF' metric. 320 Kurtosis can be defined as a measure of how prone a distribution is to outliers (Pearson, 1905). 321 SPAH4 offers a more accurate perspective by questioning not only the match of the histograms but also the extreme values and spread in the data. The 4<sup>th</sup> component is symbolized by the 322 323 expression Kurt and  $\kappa$  is the ratio of the kurtosis coefficients of the simulated (s) and observed 324 (o) data. Eq. (7) shows SPAtial Kurtosis Efficiency (SPAK) which is a three-component metric 325 replacing the histogram intersection component in the SPAEF metric with the kurtosis 326 coefficient component. Thus, it dominates the metric on its affinity for discrete values without 327 questioning histogram intersection.  $\alpha$  and  $\beta$  were introduced and explained in previous metrics, 328 also  $\kappa$  is declared in Eq. (7) as ratio of kurtosis coefficient. This metric can be characterized as 329 a mixture of SPAH4 and SPAEF metrics. Eq. (8) shows SPAtial Hybrid 5 Efficiency (SPAH5) which is a five-component metric adding skewness to the SPAH4 metric. Skewness can be 330 331 defined as a measure of the asymmetry of the data around the sample mean.

Eq. (9) shows SPAtial Histogram Equalization Efficiency (SPAHE) that is very similar to SPAEF with additional step before histogram match calculation "histogram equalization" approach. This approach is a computer image processing technique used to improve contrast in raster data. Its quantitative logic is based on the grayscale transformation (*T*) to minimize  $|c_1(T(k)) - c_0(k)|, c_0$  is the cumulative histogram of the input data, and  $c_1$  is the cumulative sum of target histogram for all intensities *k*. Histogram equalization is a specific case of the histogram remapping methods. It is an image processing technique used to advance contrast in images which spatial patterns for this study. It achieves this by efficaciously sprawling out the
 most frequent intensity values, i.e. expanding the intensity range of the image (Efford, 2000).

Eq. (10) shows SPAtial Efficiency Movers' Distance (SPAMD) is another SPAEF-oriented 341 342 multi-variate metric which measures the quantitative closeness of two pattern set by 343 considering the Earth Movers' Distance of their histograms (Rubner et al., 1998). The aim of 344 EMD approach is minimization of overall transfer cost in the conversion one histograms to another. In Eq (10),  $f_{i,i}$  is flow cost of transfer ith term of histogram K of observed map to jth 345 histogram L simulated map at distance  $d_{i,i}$ . EMD is the ratio of work done through the total 346 optimal flow and the total flow. The value of EMD is zero indicates the perfect consistency 347 348 between two histograms.

349 Eq. (11) shows Spatial Pattern Efficiency Metric (SPEM), a metric inspired by KGE and 350 SPAEF, is one of the existing metrics included in our analysis (Dembélé et al., 2020). It forces 351 the z-scores of simulated variables and observed variables to be equal (i.e., minimizing their 352 ERMS) corresponds to matching their grid cell locations (i.e., spatial patterns). SPEM 353 considers a modeled variable (Xmod) and an observed variable (Xobs) of n elements, it is 354 defined as Eq. (11); where rs is the Spearman rank-order correlation coefficient with d the difference between the ranks of Xmod and Xobs.  $\gamma$  is the variability ratio that assesses the 355 356 similarity in the dispersion of the probability distributions of Xmod and Xobs, with  $\mu$  and  $\sigma$ representing the mean and the standard deviation, respectively, and  $\alpha$  the spatial location 357 358 matching term calculated as the root-mean-square error (ERMS) of the standardized values (z-359 scores, ZX) of Xmod and Xobs (Dembélé et al., 2020). The formula for d can be written as  $d = diff(rank(X_s), rank(X_o))$ . SPEM ranges from  $-\infty$  to 1, which is its optimal value. 360

Lastly, Eq. (12) shows Structural Similarity index (SSIM) (Wang et al., 2004). An image 361 quality metric SSIM to evaluate degradation grade caused by visual data processing. This 362 363 method considers pattern similarity as it detects changes in the variation of structural 364 information between the two images. The algorithm formulates perception sensibility to visual changes based on the distortion luminance, contrast and structure information. By combining 365 three components, similarity can be characterized with overall unit metric in terms of statistical 366 properties of simulated and observed data such as mean  $\mu$ , standard deviation  $\sigma$  and covariance 367  $cov_{o,s}$ , as shown in Eq. (12).  $c_1$ ,  $c_2$  are constants that stabiles functions when the dominator 368 terms are close to zero. The SSIM is a fully referenced objective quality metric that gives values 369 370 in the range [0,1] relative to the structural relationship between the two images.

$$SSIM = \frac{(2\mu_o\mu_s + c_1)(2cov_{o,s} + c_2)}{(\mu_o^2 + \mu_s^2 + c_1)(\sigma_o^2 + \sigma_s^2 + c_2)}$$
(12)

371

All nine spatial metrics were calculated separately as long term (2002-2014) monthly average of AET data for three periods covering the growing season and combined as in Eq (13) to minimize the total error, representing objective function (OF). These periods are symbolized as March-April-May (MAM), June-July-August (JJA), and September-October-November (SON).

$$Minimize \left[ (1 - METRIC_{MAM})^2 + (1 - METRIC_{JJA})^2 + (1 - METRIC_{SON})^2 \right]$$
(13)

377 It should be noted that although we tested other metrics and approaches, we only reported nine 378 selected metrics in this study. For instance, we used harmonic mean or geometric mean instead 379 of the arithmetic mean in the second component of SPAEF. In another attempt, we replaced 380 the skewness coefficient ratio with different L-moments. We also used Hausdorff distance 381 (Hausdorff, 1914) and Fréchet distance (Fréchet, 1906) as third component in SPAEF. Even we used the product of components i.e. multiplied them instead of adding them. However, all 382 383 these attempts did not reveal better results than those reported in this study. Therefore, for 384 brevity we reported the ranking of only these nine metrics above. In this calibration study, we fine-tuned only 20 parameters of daily mHM for the Mosel Basin using the popular global 385 386 search algorithm Pareto-Archived Dynamically Dimensioned Search (ParaPADDS) algorithm (Asadzadeh & Tolson, 2013) using 750 maximum iteration and 3 parallel cores. The 20 387 388 parameters out of 69 mHM parameters are selected based on a sensitivity analysis done in our 389 previous study. Note that ParaPADDS is the multi-objective version of the Dynamically 390 Dimension Search algorithm (Tolson & Shoemaker, 2007) available in OSTRICH 391 Optimization Software Toolkit (L.S. Matott, 2017).

### 392 **5 Results**

393 In this study, six novel metrics are proposed and compared with existing SPAEF, SPEM and 394 SSIM metrics in pattern analysis of distributed hydrologic model simulations. The new metrics 395 can be called as "the sisters of SPAEF" as they have emerged from the well-established SPAEF 396 with additional unique statistical features such as automated number of bins, kurtosis and 397 skewness included in their structure. We ranked the nine metrics based on their effectiveness 398 in distinguishing between two raster maps during distributed model calibration with MODIS-399 LAI and TSEB AET for a period of 13 years from 2002 to 2014. Pre-selected 20 mHM 400 parameters are included in the following three different pattern-only calibration cases: (1) 100 401 iterations with satellite data, (2) 1000 iterations with satellite data, and (3) 1000 iterations with 402 synthetic maps. Synthetic map represents a map simulated with a known mHM parameter set for a randomly selected day that is used as the target in parameter optimization (calibration) 403 process. The use of this synthetic scenario is planned to ensure the reproducibility of the 404 405 analysis and to have a fully controlled numerical experiment. Obviously, long term monthly 406 averaging was done only with real satellite data to form robust seasonal pattern maps i.e. target 407 in the calibration.

408 Although water balance metrics, i.e. temporal metrics, are not included in the calibration, KGE 409 values are calculated to evaluate the model simulations together with standard SPAEF in Table 410 3. Streamflow simulation performance was calculated for the calibration period (2002-2014), 411 using the KGE metric between the observed gauge streamflow and simulated streamflow from 412 the model. This is done only for case 1 (TSEB 100 runs) and 2 (TSEB 1000runs) i.e. real satellite data are used in the pattern-based optimization. It is interesting to note that some of 413 414 the pattern metrics help to improve the bias in water balance as well. The three OF columns in 415 this table show lowest (best) values of each metrics reached using Eq. (13). This is particularly 416 important to show the skill of the nine metrics in converging to zero i.e. certainly exists in the synthetic case (3). It should be noted that the metrics are ranked based on the standard SPAEF 417 418 values. Closer inspection of the Table 3 shows that TSEB 1000 iterations significantly 419 improves the SPAH4 performance from 0.608 to 0.688 (SPAEF value) as compared to the TSEB 100 iterations. The reduction in OF is even more remarkable since the error in SPAH4 420 421 was halved from 0.70 to 0.35 when iterations are increased to 1000. It is clear from this table 422 that SPAHE and SPAH5 are the worst performing two metrics among all three cases.

423 Comparing the two results (100 runs vs 1000 runs) it can be seen that all metrics are improved 424 with the increased number of iterations showing the importance of the selecting appropriate 425 number of the iterations for the search algorithm. However, if enough freedom is not given to 426 the optimizer, it may fail to find the global optimum point in the solution space. Combining

427 kurtosis with skewness in the same metric (SPAH5) did not produce a discriminative metric.

428 This result is somewhat counterintuitive as we expect more constrain would yield improved

429 performance. What is striking about the values in this table is histogram equalization step did

430 not help to improve the pattern results and discriminative power of the metric.

431 **Table 3.** Calibration results of the three cases. Note that metrics are ranked based on 1000 run
432 - SPAEF values (4<sup>th</sup> numeric column).

Metrics	TSEB 100 runs			TSEF	8 1000 run	S	SYNTHETIC MAP 1000 runs		
	SPAEF	KGE	OF	SPAEF	KGE	OF	SPAEF	OF	
SPAH4	0.608	0.78	0.70	0.688	0.77	0.35	0.948	0.05	
SPACD	0.619	0.26	0.40	0.673	0.74	0.27	0.939	0.04	
SPAEF'	0.585	0.36	0.52	0.671	0.52	0.33	0.949	0.05	
SPAK	0.558	0.89	0.39	0.638	0.87	0.25	0.906	0.01	
SPAMD	0.614	0.07	0.29	0.625	0.66	0.21	0.859	0.02	
SSIM	0.557	0.21	0.19	0.491	0.41	0.15	0.948	0.00	
SPEM	0.609	0.33	1,71	0.460	0.61	1,46	0.941	0.05	
SPAHE	0.492	0.70	0.25	0.376	0.65	0.21	0.758	0.04	
SPAH5	-0.519	0.61	8,15	0.211	0.53	2,07	0.953	0.05	

### 433

434 What stands out in the table is that SSIM seems to be the most tolerant metric reaching lowest 435 OF values which corresponds to the poor SPAEF performance in all three cases. In case 3, in 436 particular, the search algorithm could converge nearly to zero SSIM but the evaluation of the 437 maps with SPAEF revealed that it is only a match around 0.95 SPAEF and not very close to 1 SPAEF i.e. perfect pattern match. In other words, minimizing SSIM in Eq (13) nearly to zero 438 439 after calibration doesn't guarantee a perfect pattern match in terms of SPAEF currency 440 (metric). Based on the results of case 1 and 2, SPAH4 and SPAK are the most successful spatial metrics for water balance. Obviously, SSIM and SPAMD have the worst KGE performance in 441 442 case 1 and 2. Note that KGE is not calculated for the synthetic case 3. Interestingly, the 443 minimization of SPEM and SPAH5 metrics via Eq (13) after optimization resulted in poor 444 values above 1 both in case 1 and 2.

445 Figure 2 shows the reference AET maps and simulated AET maps from the mHM with 446 calibrated parameters after 100 iterations (case 1). The reference three maps are given in both 447 columns for ease of comparison. The order of the metrics is in accordance with the performance 448 ranking in Table 3 and also, the ranking is provided (e.g. #1, #2 etc.) to help to the reader. The 449 combined SPAEF values of three periods (MAM, JJA and SON) are presented in brackets 450 underneath the metric name. To use a single legend, the maps are normalized with their mean. 451 The resultant maps from SPACD and SPAMD (second row in Figure 2) are slightly better than other rows as visually more similar to the reference maps (first row in Figure 2). Closer 452 453 inspection of the maps shows that the high contrast between west and south of the basin in SON period is well-captured by most of the metrics except for the SPAH5 (row 6, rank #9). 454



458 Figure 2. Long term average three-monthly TSEB reference maps versus mHM simulated 459 maps using MODIS-LAI and best-balanced Pareto solution parameter set from 100 run case.

460 Figure 3 shows the reference AET maps and simulated AET maps from the mHM with 461 calibrated parameters after 1000 iterations (case 2). It is consistent with Figure 2 that the 462 simulated AET maps by the model parameter sets optimized with SPAH4 and SPACD metrics 463 are most close to the reference maps. Similarly, the poor AET performance of SPAH5 maps is 464 apparent from the maps in the last row of the figure. Map illustration of each period reveals 465 that the combined metric value (OF) can hinder individual map performance. For instance, the SON map of the SPAHE metric in Figure 3 shows that the model better converges to the 466 remotely sensed reference map when optimized with SPAHE whereas the MAM and JJA maps 467 468 show that the model could not reproduce the AET maps of these periods as successful as with 469 the other metrics.



472 Figure 3. Long term average three-monthly TSEB reference maps versus mHM simulated 473 maps using best-balanced Pareto solution parameter set from 1000 run case.

474 The entire calibration development process, the model improvements from beginning to end 475 and the optimum points are depicted with scatter diagrams in Figure 4. It shows the relationship 476 between the value and iteration based on the ParaPADDS search algorithm, more specifically, 477 the objective function value achieved for each iteration step of the calibration process. While 478 the OF results in Table 3 are obtained at the end of the iteration step sequence, some consistent 479 metrics may reach this best value earlier. SPAH4 reached its best OF value at 0.70 and 0.35 in 480 approximately quarter steps for 100 and 1000 runs, respectively. Similarly, SPACD, SPAEF' 481 and SPAMD are also fast-improving metrics. Since the synthetic case was based on a virtually 482 generated daily map, it took longer for the metrics to find the points where their improvement 483 became linear, nearly a third. It is surprising to see that SPAH5 and SPEM are consistent early 484 maturing metrics despite their poor spatial performance.







Figure 4. Scatter plots of the calibration processes, the OF value-iteration relationship of the 488 489 PDSS search algorithm. First and second column sub-plots are the same figures except for 490 different extent.



492 Figure 5. Monthly average hydrograph of all years in the calibration period (2002–2014) to
 493 demonstrate the flow simulation performances of nine different metrics.

494 Figure 5 compares in-situ observed hydrographs and simulated hydrographs constrained by 495 metrics. SPAH4 and SPAK performed better in each case, predicting the most similar 496 discharges to the observed Cochem outflows. Otherwise, the SPAHE metric standout for the 497 100 runs and the SPACD metric for the 1000 runs, as pointed out by the KGE column in Table 498 3. The simulations show better hydrograph fitting during the growing season, especially during 499 the summer months, also the hydrograph line breakpoints, peaks and valleys are coherently 500 followed. Thus, the overall trend and characteristics of the streamflow were successfully 501 analyzed and represented. Also, a positive correlation was found between increasing iteration 502 and hydrograph fit. As the number of iterations increases, the hydrograph lines become closer to the observed lines and the overall consolidation of the hydrographs provides better results. 503 504 The narrow range of hydrographs in Figure 5 shows that the developed new metrics can be 505 used not only for the spatial pattern performance simulating the AET but also for the temporal 506 streamflow performance simulating the discharges.

507 Overall, the results indicate that the newly developed SPAH4 and SPACD are the best 508 performing metrics for all calibration scenarios, particularly in the non-synthetic TSEB cases. 509 The competitive performance of the SPAMD metric that follows them should not be ignored. 510 Briefly, the four-component spatial performance metric SPAH4 stands out especially with its 511 versatile evaluation and robust performance, indicated with bold text in Table 3. Although the 512 modeler can use the SPAH4 and SPACD metrics in the long and short runs, respectively, both 513 offer close values for the decision makers. We can see that the only negative output is experienced in the TSEB 100 runs i.e. SPAH5 It should not be overlooked that SPAH5 is a 514 515 prominent metric for synthetic scenarios. Interestingly, there is a significant positive correlation 516 between the KGE and the metrics containing the kurtosis statistic.

### 517 6 Discussion

518 This study sets out to assess the importance and comparison of spatial metrics in distributed 519 model calibration. Previous studies have noted that spatial metrics are closer to the reference model than time series metrics in model optimizations (Demirel et al., 2018). One of the first 520 521 objectives of the study is to select the appropriate spatial performance metric that plays an 522 active role in simulating inadequate spatial AET models similar to satellite-based reference 523 models. SPAEF has been the inspiration for this study with its innovations in spatial model parameterization and spatial performance metric selection. These innovations have raised new 524 525 questions in the pattern comparison used in model optimization. Numerous imperfect models 526 are produced during these optimizations, due to limitations in the chosen objective function. 527 To overcome these limitations and to obtain a more physically meaningful and empowered 528 metric, we have developed new metrics that include statistical and analytical approaches. 529 Thanks to this meta-analysis, while suggesting the most successful metric for users, different 530 objective functions that can be used for various purposes can also be seen as an opportunity. 531 While searching for new solutions for a more robust spatial performance metric, we derived 532 metrics that emphasize spatiality in a more comprehensive way by increasing the number of 533 components of SPAEF and changing the content of the components. For the three cases, 534 significant findings that are both different from each other and support each other have been 535 identified. The TSEB 100 and 1000 run cases in model calibration served the purpose of 536 evaluating metric performances in short and long runs, thus providing a flexible and versatile 537 assessment that allows the progress of the model calibration performed by the metrics to be 538 monitored and the decision maker to choose metrics according to their preferences.

539 TSEB 100 runs, which we tested by focusing on the performance of spatial metrics in short 540 runs, SPACD and SPAMD demonstrated better results on the SPAEF basis compared to other 541 metrics. Notably, SPAK and SPAH4 including the kurtosis coefficient ratio component, 542 yielded the best KGE values even at iterations close to the beginning. TSEB 1000 runs which 543 we tested by focusing on its performance in long runs, resulted in more decisive outcomes with 544 no negative values for any criteria. SPAH4 emerged as the top-performing metric in this case, 545 followed by SPACD. The competition between these metrics was notable. In the uncertainty 546 analysis, SPAH4 has an acceptable sampling error although it has the extra component. (Table A1). Like the TSEB 100 runs, SPAK and SPAH4 exhibited the highest KGE values. This 547 548 indicated consistency was strong evidence for important findings and suggests that the 549 descriptive statistical kurtosis ratio component has a considerable positive effect on the discharge simulation. Due to the tendency of the SPAH4 metric including kurtosis for flow 550 551 prediction, it worked as a metric that focused on both spatial and flow performance, although 552 the analysis was performed with a single spatial performance-oriented objective function. It 553 sheds light on the analysis in detecting the presence of outliers potential also differences in the tail and crests, controlling data integrity, understanding data distribution, reliability of the 554 555 statistical analysis and improving the metric performance from a statistical perspective. Thus, by investigating and questioning the effect of outliers on spatial performance, the harmony and 556 557 differences between them are also included in the model. Now that these outliers are introduced to the model, the histogram intercept component is also supported, the margin of error is 558 559 reduced and a more exact match is made.

560 In the synthetic scenario, the metric SPAH5 which incorporates skewness characteristics, 561 yielded the best SPAEF value. SPACD and SPAH4 also demonstrated successful outcomes in 562 this scenario. The kurtosis information we use in the SPAH4 metric expresses how often 563 outliers occur, while the skewness information we use as the fifth component in the SPAH5 564 metric gives information about the direction of the outliers. Our purpose in including the 565 skewness component is to question the likelihood of events in the probability distribution, and 566 especially to consider extreme distribution. Various datasets have different characteristics, since the differences specific to this dataset represent important concepts in the calibration 567 568 model, many principles are referred to using the skewness information, from the algorithm of the model to the physics-based hydrology information. Thus, we enabled a more 569 570 comprehensive and more specific analysis for models consisting of diverse data. Our finding 571 of the importance of these statistical measures in understanding the data is supported by the 572 study by Cain et al., processing skewness and kurtosis information on distributions collected 573 from the authors of the published articles (Cain et al., 2017). In addition, it is possible to derive 574 a positive interpretation from a negative finding in meta-analyses as in this study. Since the 575 only difference between the metrics with the best and the worst performance in TSEB runs, 576 namely SPAH4 and SPAH5, is the skewness ratio component, it can be concluded that 577 skewness is a component that negatively affects the spatial metrics used in pattern comparison. 578 It should be noted though that skewness information is an outstanding component for synthetic 579 cases.

580 In TSEB 100 runs scenario, the spatial performance tussle results of the metrics show that the 581 newly proposed metric i.e. SPACD outperforms the conventional three-component metric 582 SPAEF (5.76% better) on the other hand 11.11% better than SSIM and 1.66% better than 583 SPEM. In TSEB 1000 runs results demonstrate that the newly developed four-component 584 metric i.e. SPAtial Hybrid 4 (SPAH4) slightly outperform SPAEF (2.62% better). However, 585 SPAH4 significantly outperforms the other existing metrics i.e. 40.22% better than SSIM and 586 49.53% better than SPEM.

### 587 7 Conclusion

588 In this study, we thoroughly assessed common existing metrics and new spatial pattern-oriented performance metrics that we developed based on SPAEF. For the consistency and reliability 589 590 of the results, the Mosel Basin with high data quality was selected and the physics-based fully 591 distributed mHM model was established for this basin. In these three different scenarios, we 592 performed analyses with various (low-high) iterations for actual evapotranspiration maps (TSEB AET) and synthetic maps. The most popular metrics (SPAEF, SSIM and SPEM) were 593 594 compared with new metrics (SPAH4, SPACD, etc.) to measure the convergence of the mHM 595 model to long-term monthly AET maps observed during parameter calibration. The usage of 596 this synthetic scenario is important to ensure the reproducibility of the experiments and to give 597 us full control over the calibration process. Based on our findings we can draw the following 598 conclusions.

The inclusion of kurtosis ratio coefficient in the spatial pattern-oriented metrics demonstrates
 that metric performance is improved, so it has a positive impact on the spatially objective
 functions. Also shows a positive effect on streamflow prediction, it successfully calibrates the
 KGE metric even in very short runs. Furthermore, while using the skewness ratio coefficient
 gave unsuccessful results for TSEB AET maps, the kurtosis information of the distribution was
 more prominent in the pattern performance of the models. However, the SPAH5 performs the
 best among the close results and is presented as a strong hypothesis for the synthetic cases.

- The metric with the best performance in the short runs was SPACD, which normalizes the
distribution according to density. The excellent consistency between histograms, which is the
main component of the Earth mover's distance metric, has a positive effect on making this
metric a sharp metric with little tolerance, making SPAMD the second-best metric.

- 610 The best-performing metric on long runs was SPAH4, a four-component spatial performance
- 611 metric that includes the kurtosis of the distribution. It was followed by the SPACD metric,
- 612 which proved its consistent performance. Thus, the decision maker is presented with a flexible
- 613 and wide working area.
- Considering all the experimental results, the most successful and robust metric in all three
   scenarios is our newly developed spatial pattern-oriented SPAH4, which outperforms the
   existing metrics in the literature by up to fifty per cent.
- 617 In future studies, it would significantly enhance the depth and quality of the analysis to increase
- 618 the number of iterations. In fact, convergence in hydrological models is closely related to the
- 619 number of parameters and the freedom of the appropriate iteration chosen. Future work may
- benefit from exploring untested statistical terms to add a new perspective. We expect that these
- 621 newly developed metrics, especially SPAH4, will be used not only in hydrology but also in
- 622 other fields including remote sensing, image processing and object detection.

623 **Appendix A:** Results of the jackknife and bootstrap based sampling uncertainty analysis. Clark et al (2021) showed that the two most popular metrics in hydrology, i.e. NSE and KGE, are 624 625 vulnerable to sampling uncertainty since the differences between observed and simulated 626 streamflow values at random time steps in time series which can have significant effects on the results (Knoben & Spieler, 2022). From this study, we are inspired to assess the sampling 627 uncertainty in ten metrics using the gumboot R package (Clark et al., 2021) which uses a 628 629 jackknife-after-bootstrap method of Efron (1992) to estimate standard errors (SEJaB) shown 630 in Table A1.

GOF_stat	seJack	seBoot	p05	p50	թ95	score	biasJack	biasBoot	seJab
SSIM	0.0103	0.0099	0.6144	0.6311	0.6457	0.6307	-0.0002	0.0000	0.0091
SPAHE	0.0568	0.0119	0.7717	0.7917	0.8107	0.7783	0.1496	0.0131	0.0112
SPAMD	0.0114	0.0108	0.6785	0.6972	0.7137	0.6966	0.0006	0.0001	0.0115
SPEM	0.0180	0.0175	0.2739	0.3041	0.3309	0.3034	-0.0006	-0.0003	0.0146
SPAEF	0.0133	0.0128	0.6489	0.6711	0.6917	0.6727	0.0017	-0.0021	0.0148
SPAEF'	0.0133	0.0127	0.6489	0.6711	0.6917	0.6727	0.0017	-0.0021	0.0152
SPAK	0.0302	0.0288	0.5719	0.6226	0.6661	0.6207	-0.0004	0.0007	0.0295
SPAH4	0.0302	0.0298	0.5484	0.5999	0.6459	0.6000	0.0011	-0.0012	0.0313
SPACD	0.0234	0.0248	0.6056	0.6571	0.6851	0.6670	-0.0219	-0.0142	0.0603
SPAH5	0.1685	0.2077	-0.3594	0.0373	0.2636	0.0427	-0.0388	-0.0401	0.3382

631 **Table A1.** Sampling uncertainty of the metrics i.e. ranked based on the seJab column.

632



Figure A1. eCDF plot of daily discharge for all years in the calibration period (2002-2014) to
 visualize the distribution of the data and identify statistical patterns.

637 Figure A1 visualizes the empirical cumulative distribution function (eCDF) plot for the 638 observed and simulated data, which shows how the probability of a given discharge value 639 occurring varies over the range of discharge values. In this context, the percentage of observed 640 discharges less than nearly 500 is 80% and less than 200 is 50% for both the TSEB 100 and 641 1000 runs. Furthermore, the slope of the curve at any point represents the density function of 642 the discharge values at that point, and the intervals where the curve steepens contain values 643 close to the mean value. Hence, it can be concluded that the overall average discharge value of the steepening intervals of the flow data resulting from the simulation of the metrics is roughly 644 645  $300 \text{ m}^3$ /s. The mean observed outflow of Cochem station is around  $315 \text{ m}^3$ /s supports this 646 outcome. In both cases, SPAK and SPAH4 illustrated a high level of matching in terms of the 647 fit of the curves generated by the observed data (OBS) and the metrics, with the least difference 648 between the distributions.





Figure A2. Monthly average hydrograph of the last two years in the calibration period (2013–
2014)

### 654 Acknowledgements, Samples, and Data

655 Data Availability Statement: Discharge data is provided by GRDC data portal (https://portal.grdc.bafg.de/) in Koblenz, Germany. MODIS MOD16A2 v061 product was 656 retrieved from https://doi.org/10.5067/MODIS/MOD16A2.061. SRTM DEM data was 657 retrieved from https://www.earthdata.nasa.gov. The source code of the mHM is publicly 658 659 available at https://doi.org/10.5281/zenodo.4575390. The source code of the SPAEF metric is publicly available at https://doi.org/10.5281/zenodo.5861253. The model calibration software 660 Ostrich is available from https://github.com/usbr/ostrich. The simulation scripts and results of 661 662 the mHM model simulations are publicly available at https://doi.org/10.5281/zenodo.8059198. The source code to quantify the sampling uncertainty in performance metrics (the "gumboot" 663 package) is available at https://github.com/CH-Earth/gumboot. The scripts and results of the 664 665 gumboot-based sampling uncertainty analysis is available at https://doi.org/10.5281/zenodo.8058659 666

Acknowledgements: We acknowledge the financial support for the SPACE project by the
Villum Foundation (http://villumfonden.dk/) through their Young Investigator Program (grant
VKR023443). The first author is supported by NASA program i.e. NNH22ZDA001N-RRNES:
A.24 Rapid Response and Novel Research in Earth Science under the grant number 22RRNES22-0010 and by the Scientific Research Projects Department of Istanbul Technical
University (ITU-BAP) under grant number MDA-2022-43762 and by the National Center for

- High Performance Computing of Turkey (UHeM) under grant number 1007292019.
- 674 **Conflicts of Interest:** "The authors declare no conflict of interest."
- 675 Institutional Review Board Statement: Not applicable.
- 676 **Informed Consent Statement:** Not applicable.
- 677
- 678

### 679 **References**

- Ahmed, K., Sachindra, D. A., Shahid, S., Demirel, M. C., & Chung, E.-S. (2019). Selection of
  multi-model ensemble of general circulation models for the simulation of precipitation
  and maximum and minimum temperature based on spatial assessment metrics. *Hydrology and Earth System Sciences*, 23(11), 4803–4824. https://doi.org/10.5194/hess-23-48032019
- Allen, R. G., Pereira, L. S., Raes, D., & Smith, M. (1998). Crop evapotranspiration -*Guidelines for computing crop water requirements*. FAO Irrigation and drainage paper
  56. Retrieved from http://www.fao.org/docrep/x0490e/x0490e00.htm (accessed at
  16/02/2018)
- Arun, N., Gaw, N., Singh, P., Chang, K., Aggarwal, M., Chen, B., et al. (2021). Assessing the
   Trustworthiness of Saliency Maps for Localizing Abnormalities in Medical Imaging.
   *Radiology: Artificial Intelligence*, 3(6). https://doi.org/10.1148/ryai.2021200267
- Asadzadeh, M., & Tolson, B. (2013). Pareto archived dynamically dimensioned search with
   hypervolume-based selection for multi-objective optimization. *Engineering Optimization*,
   45(12), 1489–1509. https://doi.org/10.1080/0305215X.2012.748046
- Avcuoğlu, M. B., & Demirel, M. C. (2022). Hidrolojik Model Kalibrasyonunda Uydu Tabanlı
  Aylık Buharlaşma ve LAI Verilerinin Kullanılması. *Teknik Dergi*, *33*(6), 13013–13035.
  https://doi.org/10.18400/tekderg.1067466
- Becker, R., Koppa, A., Schulz, S., Usman, M., aus der Beek, T., & Schüth, C. (2019). Spatially
  distributed model calibration of a highly managed hydrological system using remote
  sensing-derived ET data. *Journal of Hydrology*, 577(10), 123944.
  https://doi.org/10.1016/j.jhydrol.2019.123944
- Beven, K. (2023). Benchmarking hydrological models for an uncertain future. *Hydrological Processes*, *37*(5). https://doi.org/10.1002/hyp.14882
- de Boer-Euser, T., Bouaziz, L., De Niel, J., Brauer, C., Dewals, B., Drogue, G., et al. (2017).
  Looking beyond general metrics for model comparison lessons from an international
  model intercomparison study. *Hydrology and Earth System Sciences*, 21(1), 423–440.
  https://doi.org/10.5194/hess-21-423-2017
- Busari, I. O., Demirel, M. C., & Newton, A. (2021). Effect of Using Multi-Year Land Use Land
  Cover and Monthly LAI Inputs on the Calibration of a Distributed Hydrologic Model. *Water*, *13*(11), 1538. https://doi.org/10.3390/w13111538
- Cain, M. K., Zhang, Z., & Yuan, K.-H. (2017). Univariate and multivariate skewness and
   kurtosis for measuring nonnormality: Prevalence, influence and estimation. *Behavior Research Methods*, 49(5), 1716–1735. https://doi.org/10.3758/s13428-016-0814-1
- Clark, M. P., Vogel, R. M., Lamontagne, J. R., Mizukami, N., Knoben, W. J. M., Tang, G., et
  al. (2021). The Abuse of Popular Performance Metrics in Hydrologic Modeling. *Water Resources Research*, 57(9). https://doi.org/10.1029/2020WR029001
- Danapour, M., Fienen, M. N., Højberg, A. L., Jensen, K. H., & Stisen, S. (2021). MultiConstrained Catchment Scale Optimization of Groundwater Abstraction Using Linear
  Programming. *Groundwater*, 59(4), 503–516. https://doi.org/10.1111/gwat.13083
- Dembélé, M., Hrachowitz, M., Savenije, H. H. G., Mariéthoz, G., & Schaefli, B. (2020).
  Improving the Predictive Skill of a Distributed Hydrological Model by Calibration on

- Spatial Patterns With Multiple Satellite Data Sets. *Water Resources Research*, 56(1).
  https://doi.org/10.1029/2019WR026085
- Demirel, M. C., Booij, M. J., & Hoekstra, A. Y. (2013). Effect of different uncertainty sources
   on the skill of 10 day ensemble low flow forecasts for two hydrological models. *Water Resources Research*, 49(7), 4035–4053. https://doi.org/10.1002/wrcr.20294
- Demirel, M. C., Booij, M. J., & Hoekstra, A. Y. (2015). The skill of seasonal ensemble low flow forecasts in the Moselle River for three different hydrological models. *Hydrology and Earth System Sciences*, 19(1), 275–291. https://doi.org/10.5194/hess-19-275-2015
- Demirel, M. C., Mai, J., Mendiguren, G., Koch, J., Samaniego, L., & Stisen, S. (2018).
  Combining satellite data and appropriate objective functions for improved spatial pattern
  performance of a distributed hydrologic model. *Hydrology and Earth System Sciences*,
  22(2), 1299–1315. https://doi.org/10.5194/hess-22-1299-2018
- Dougherty, E., Sherman, E., & Rasmussen, K. L. (2020). Future Changes in the Hydrologic
  Cycle Associated with Flood-Producing Storms in California. *Journal of Hydrometeorology*, 21(11), 2607–2621. https://doi.org/10.1175/JHM-D-20-0067.1
- Efford, N. (2000). Digital Image Processing: A Practical Introduction Using JavaTM. Pearson
   Education. Slate.
- Efron, B. (1992). Jackknife-After-Bootstrap Standard Errors and Influence Functions. *Journal*of the Royal Statistical Society: Series B (Methodological), 54(1), 83–111.
  https://doi.org/10.1111/j.2517-6161.1992.tb01866.x
- Fréchet, M. M. (1906). Sur quelques points du calcul fonctionnel. *Rendiconti Del Circolo Matematico Di Palermo*, 22(1), 1–72. https://doi.org/10.1007/BF03018603
- Freedman, D., & Diaconis, P. (1981). On the histogram as a density estimator:L 2 theory. *Zeitschrift Für Wahrscheinlichkeitstheorie Und Verwandte Gebiete*, 57(4), 453–476.
  https://doi.org/10.1007/BF01025868
- Gaur, S., Singh, B., Bandyopadhyay, A., Stisen, S., & Singh, R. (2022). Spatial pattern-based
   performance evaluation and uncertainty analysis of a distributed hydrological model.
   *Hydrological Processes*, *36*(5). https://doi.org/10.1002/hyp.14586
- Gupta, H. V., Kling, H., Yilmaz, K. K., & Martinez, G. F. (2009). Decomposition of the mean
  squared error and NSE performance criteria: Implications for improving hydrological
  modelling. *Journal of Hydrology*, 377(1–2), 80–91.
  https://doi.org/10.1016/j.jhydrol.2009.08.003
- Hargreaves, G. H., & Samani, Z. A. (1985). Reference Crop Evapotranspiration from
  Temperature. *Applied Engineering in Agriculture*, 1(2), 96–99.
  https://doi.org/10.13031/2013.26773
- Hausdorff, F. (1914). Bemerkung über den Inhalt von Punktmengen. *Mathematische Annalen*,
  758 75(3), 428–433.
- Hossain, M. K., & Meng, Q. (2020). A thematic mapping method to assess and analyze
  potential urban hazards and risks caused by flooding. *Computers, Environment and Urban Systems*, 79, 101417. https://doi.org/10.1016/j.compenvurbsys.2019.101417
- Immerzeel, W. W., & Droogers, P. (2008). Calibration of a distributed hydrological model
  based on satellite evapotranspiration. *Journal of Hydrology*, 349(3–4), 411–424.

- 764 https://doi.org/10.1016/j.jhydrol.2007.11.017
- Knoben, W. J. M., & Spieler, D. (2022). Teaching hydrological modelling: illustrating model
  structure uncertainty with a ready-to-use computational exercise. *Hydrology and Earth System Sciences*, 26(12), 3299–3314. https://doi.org/10.5194/hess-26-3299-2022
- Knoben, W. J. M., Freer, J. E., & Woods, R. A. (2019). Technical note: Inherent benchmark or
   not? Comparing Nash-Sutcliffe and Kling-Gupta efficiency scores. *Hydrology and Earth System Sciences Discussions*, (July), 1–7. https://doi.org/10.5194/hess-2019-327
- Kumar, R., Samaniego, L., & Attinger, S. (2013). Implications of distributed hydrologic model
   parameterization on water fluxes at multiple scales and locations. *Water Resources Research*, 49(1), 360–379. https://doi.org/10.1029/2012WR012195
- López, P., Sutanudjaja, E. H., Schellekens, J., Sterk, G., & Bierkens, M. F. P. (2017).
  Calibration of a large-scale hydrological model using satellite-based soil moisture and
  evapotranspiration products. *Hydrology and Earth System Sciences*, 21(6), 3125–3144.
  https://doi.org/10.5194/hess-21-3125-2017
- Martinez-Villalobos, C., Neelin, J. D., & Pendergrass, A. G. (2022). Metrics for Evaluating
   CMIP6 Representation of Daily Precipitation Probability Distributions. *Journal of Climate*, *35*(17), 5719–5743. https://doi.org/10.1175/JCLI-D-21-0617.1
- Matott, L. Shawn. (2004). OSTRICH: an Optimization Software Tool, Documentation and
   User's Guide, Version 17.12.19. Retrieved from https://github.com/usbr/ostrich
- Matott, L.S. (2017). OSTRICH: an Optimization Software Tool, Documentation and User's
   Guide. University at Buffalo Center for Computational Research, Version 17, 79.
- Monteith, J. L. (1965). EVAPORATION AND ENVIRONMENT. Symposia of the Society for
   *Experimental Biology*, 19, 205–234.
- Nash, J. E., & Sutcliffe, J. V. (1970). River flow forecasting through conceptual models part I
  A discussion of principles. *Journal of Hydrology*, 10(3), 282–290.
  https://doi.org/10.1016/0022-1694(70)90255-6
- Nilsson, J., & Akenine-Möller, T. (2020). Understanding SSIM. Retrieved from https://arxiv.org/abs/2006.13846
- Norman, J. M., Kustas, W. P., & Humes, K. S. (1995). Source approach for estimating soil and
   vegetation energy fluxes in observations of directional radiometric surface temperature.
   *Agricultural and Forest Meteorology*, 77(3–4), 263–293. https://doi.org/10.1016/0168 1923(95)02265-Y
- Odusanya, A. E., Schulz, K., & Mehdi-Schulz, B. (2022). Using a regionalisation approach to
  evaluate streamflow simulated by an ecohydrological model calibrated with global land
  surface evaporation from remote sensing. *Journal of Hydrology: Regional Studies*, 40,
  101042. https://doi.org/10.1016/j.ejrh.2022.101042
- 800 Onyutha, C. (2022). A hydrological model skill score and revised R-squared. *Hydrology* 801 *Research*, 53(1), 51–64. https://doi.org/10.2166/nh.2021.071
- 802 Pearson, K. (1905). The problem of the random walk. *Nature*, 72(1865), 294.
- Priestley, C. H. B., & Taylor, R. J. (1972). On the Assessment of Surface Heat Flux and
   Evaporation Using Large-Scale Parameters. *Monthly Weather Review*, 100(2), 81–92.

- 805 https://doi.org/10.1175/1520-0493(1972)100<0081:OTAOSH>2.3.CO;2
- Rakovec, O., Kumar, R., Mai, J., Cuntz, M., Thober, S., Zink, M., et al. (2016). Multiscale and
   Multivariate Evaluation of Water Fluxes and States over European River Basins. *Journal of Hydrometeorology*, *17*(1), 287–307. https://doi.org/10.1175/JHM-D-15-0054.1
- Rientjes, T. H. M., Muthuwatta, L. P., Bos, M. G., Booij, M. J., & Bhatti, H. A. (2013). Multi-variable calibration of a semi-distributed hydrological model using streamflow data and satellite-based evapotranspiration. *Journal of Hydrology*, 505, 276–290. https://doi.org/10.1016/j.jhydrol.2013.10.006
- Rubner, Y., Tomasi, C., & Guibas, L. J. (1998). A metric for distributions with applications to
  image databases. In *Sixth International Conference on Computer Vision (IEEE Cat. No.98CH36271)* (pp. 59–66). Narosa Publishing House.
  https://doi.org/10.1109/ICCV.1998.710701
- Samaniego, L., Kumar, R., & Attinger, S. (2010). Multiscale parameter regionalization of a
  grid-based hydrologic model at the mesoscale. *Water Resources Research*, 46(5),
  W05523. https://doi.org/10.1029/2008WR007327
- Samaniego, L., Brenner, J., Craven, J., Cuntz, M., Dalmasso, G., Demirel, C. M., et al. (2021,
  July 21). The mesoscale Hydrologic Model mHM v5.11.2.
  https://doi.org/10.5281/ZENODO.5119952
- Schneider, R., Henriksen, H. J., & Stisen, S. (2022). A robust objective function for calibration
  of groundwater models in light of deficiencies of model structure and observations. *Journal of Hydrology*, *613*, 128339. https://doi.org/10.1016/j.jhydrol.2022.128339
- Scott, D. W. (1979). On optimal and data-based histograms. *Biometrika*, 66(3), 605–610.
  https://doi.org/10.1093/biomet/66.3.605
- Sirisena, T. A. J. G., Maskey, S., & Ranasinghe, R. (2020). Hydrological Model Calibration
  with Streamflow and Remote Sensing Based Evapotranspiration Data in a Data Poor
  Basin. *Remote Sensing*, *12*(22), 3768. https://doi.org/10.3390/rs12223768
- Stisen, S., Koch, J., Sonnenborg, T. O., Refsgaard, J. C., Bircher, S., Ringgaard, R., & Jensen,
  K. H. (2018). Moving beyond run-off calibration-Multivariable optimization of a surfacesubsurface-atmosphere model. *Hydrological Processes*, 32(17), 2654–2668.
  https://doi.org/10.1002/hyp.13177
- Sturges, H. A. (1926). The Choice of a Class Interval. *Journal of the American Statistical Association*, 21(153), 65–66. https://doi.org/10.1080/01621459.1926.10502161
- Swain, M. J., & Ballard, D. H. (1991). Color indexing. *International Journal of Computer Vision*, 7(1), 11–32. https://doi.org/10.1007/BF00130487
- Thober, S., Cuntz, M., Kelbling, M., Kumar, R., Mai, J., & Samaniego, L. (2019). The
  multiscale routing model mRM v1.0: simple river routing at resolutions from 1 to 50 km. *Geoscientific Model Development*, 12(6), 2501–2521. https://doi.org/10.5194/gmd-122501-2019
- Thoya, P., Maina, J., Möllmann, C., & Schiele, K. S. (2021). AIS and VMS Ensemble Can
  Address Data Gaps on Fisheries for Marine Spatial Planning. *Sustainability*, *13*(7), 3769.
  https://doi.org/10.3390/su13073769
- 846 Tolson, B. A., & Shoemaker, C. A. (2007). Dynamically dimensioned search algorithm for

- computationally efficient watershed model calibration. *Water Resources Research*, 43(1).
  https://doi.org/10.1029/2005WR004723
- Wakigari, S. A., & Leconte, R. (2023). Assessing the Potential of Combined SMAP and InSitu Soil Moisture for Improving Streamflow Forecast. *Hydrology*, *10*(2), 31.
  https://doi.org/10.3390/hydrology10020031
- Wang, Z., Bovik, A. C., Sheikh, H. R., & Simoncelli, E. P. (2004). Image Quality Assessment:
  From Error Visibility to Structural Similarity. *IEEE Transactions on Image Processing*, *13*(4), 600–612. https://doi.org/10.1109/TIP.2003.819861
- Wiederholt, R., Paudel, R., Khare, Y., Davis, S. E., Melodie Naja, G., Romañach, S., et al.
  (2019). A multi-indicator spatial similarity approach for evaluating ecological restoration
  scenarios. *Landscape Ecology*, *34*(11), 2557–2574. https://doi.org/10.1007/s10980-01900904-w
- Yoo, S. B. M., Tu, J. C., Piantadosi, S. T., & Hayden, B. Y. (2020). The neural basis of
   predictive pursuit. *Nature Neuroscience*, 23(2), 252–259. https://doi.org/10.1038/s41593 019-0561-6
- Zink, M., Mai, J., Cuntz, M., & Samaniego, L. (2018). Conditioning a Hydrologic Model Using
  Patterns of Remotely Sensed Land Surface Temperature. *Water Resources Research*,
  54(4), 2976–2998. https://doi.org/10.1002/2017WR021346