



Manuscript in submission to *ESSD*

1 **A Synthesis of Global Streamflow characteristics, Hydrometeorology, and**
2 **catchment Attributes (GSHA) for Large Sample River-Centric Studies**

3
4 Ziyun Yin¹, Peirong Lin^{1,2*}, Ryan Riggs³, George H. Allen⁴, Xiangyong Lei¹, Ziyuan
5 Zheng^{5,6}, Siyu Cai⁷

- 6 1. Institute of Remote Sensing and GIS, School of Earth and Space Sciences, Peking University
7 2. International Research Center for Big Data for Sustainable Development Goals, Beijing, China
8 3. Department of Geography, Texas A&M University, Texas, USA
9 4. Department of Geosciences, Virginia Polytechnic Institute and State University, Virginia, USA
10 5. Key Laboratory of Regional Climate-Environment Research for Temperate East Asia, Institute of
11 Atmospheric Physics, Chinese Academy of Sciences, Beijing, China
12 6. University of Chinese Academy of Sciences, Beijing, China
13 7. State Key Laboratory of Simulation and Regulation of Water Cycle in River Basin, China Institute
14 of Water Resources and Hydropower Research, Beijing, China

15 * *Correspondence to:* Peirong Lin (peironglinlin@pku.edu.cn)

16 Manuscript submitted to *ESSD*, July 1st, 2023

17 **Abstract**

18 Our understanding and predictive capability of streamflow processes largely rely on high-
19 quality datasets that depict a river's upstream basin characteristics. Recent proliferation of large
20 sample hydrology (LSH) datasets has promoted model parameter estimation and data-driven
21 analyses of the hydrological processes worldwide, yet existing LSH is still insufficient in terms of
22 sample coverage, uncertainty estimates, and dynamic descriptions of anthropogenic activities. To
23 bridge the gap, we contribute the Synthesis of Global Streamflow characteristics, Hydrometeorology,
24 and catchment Attributes (GSHA) to complement existing LSH datasets, which covers 21,568
25 watersheds from 13 agencies for as long as 43 years based on discharge observations scraped from
26 web. In addition to annual streamflow indices, each basin's daily meteorological variables (i.e.,
27 precipitation, 2 m air temperature, longwave/shortwave radiation, wind speed, actual and potential
28 evapotranspiration), daily-weekly water storage terms (i.e., snow water equivalence, soil moisture,
29 groundwater percentage), and yearly dynamic descriptors of the land surface characteristics (i.e.,
30 urban/cropland/forest fractions, leaf area index, reservoir storage and degree of regulation) are also
31 provided by combining openly available remote sensing and reanalysis datasets. The uncertainties
32 of all meteorological variables are estimated with independent data sources. Our analyses revealed
33 the following insights: (i) meteorological data uncertainties vary across variables and geographical
34 regions, and the prominent patterns revealed should be accounted for by LSH users, (ii) ~6%
35 watersheds shifted between human managed and natural states during the GSHA time span, which
36 may be useful for analysis that takes the changing land surface characteristics into account, and (iii)
37 GSHA watersheds observed a more widespread declining trend in runoff coefficient than an
38 increasing trend, pointing towards critical water availability issues. Overall, GSHA is expected to
39 serve hydrological model parameter estimation and data-driven analyses as it continues to improve.
40 GSHA v1.0 can be accessed at <https://doi.org/10.5281/zenodo.8090704> (Yin et al., 2023).



Manuscript in submission to *ESSD*

41 1 Introduction

42 Climate change has posed profound challenges to the management of freshwater resources,
43 specifically riverine floods and water shortages (AghaKouchak et al., 2020; Thackeray et al., 2022).
44 The urgent need for flood and drought forecasting, water resources planning and management, all
45 call for high-quality streamflow predictions for basins worldwide to analyse global terrestrial water
46 conditions in a systematic view (Burges, 1998). The scarcity of hydrological observations has
47 brought challenges to these predictions (Belvederesi et al., 2022; Hrachowitz et al., 2013), thus the
48 development of computer models that allow for “modelling everything everywhere” (Beven &
49 Alcock, 2012) constitutes the backbone of hydrological studies. Existing studies have used
50 physically-based and data-driven models for streamflow simulation (Lin et al., 2018; Nandi &
51 Reddy, 2022; Zhang et al., 2020), with efforts to improve accuracy of prediction by combining both
52 (Cho & Kim, 2022; Razavi & Coulibaly, 2013). Yet the prediction of the magnitude, timing, and
53 trend of critical streamflow characteristics are still subject to multiple sources of errors and
54 uncertainties (Bourdin et al., 2012; Brunner et al., 2021).

55 Streamflow (Q) can be represented by the simple water balance equation involving
56 precipitation (P), evapotranspiration (ET), and water storage terms (S) denoted as $Q = P - ET - \Delta S$,
57 yet influencing factors of these components could bring uncertainties that will cascade downstream.
58 Starting from the model assumptions to the data used to represent climate, soil water, ice cover,
59 topography and land use, as well as to the less well-known processes such as human perturbations
60 and sub-surface flows (Benke et al., 2008; Wilby & Dessai, 2010), these complications impede our
61 understanding of streamflow processes across scales, which also limits the modelling and predictive
62 capability for streamflow. Thus, reducing the predictive uncertainties require high-quality data with
63 massive samples capable of depicting each of the water balance components, as well as the natural
64 and anthropogenic factors involved (Gupta et al., 2014).

65 Efforts have been made to address the need of such kind of high-quality datasets on watershed-
66 scale hydro-climate and environmental conditions during the past couple of decades. One of the
67 earliest was the most widely used dataset generated for the Model Parameter Estimation Experiment
68 (MOPEX) project aimed at better hydrological modelling (Duan et al., 2006). Historical hydro-
69 meteorological data and land surface characteristics for over 400 hydrologic basins in the United
70 States were provided, which was fundamental to the progress in large sample hydrology (LSH)
71 (Addor et al., 2020; Schaake et al., 2006). Later the dataset was expanded to 671 catchments in the
72 contiguous United States (CONUS) and benchmarked by model results (Newman et al., 2015).
73 Based on these studies, the Catchment Attributes and Meteorology for Large-sample Studies
74 (CAMELS) dataset was developed, providing comprehensive and updated data on topography,
75 climate, streamflow, land cover, soil, and geology attributes for each catchment (Addor et al., 2017).
76 The CONUS CAMELS dataset soon became influential in LSH and has since inspired researchers
77 from Australia (Fowler et al., 2021), Europe (Coxon et al., 2020; Delaigue et al., 2022; Klingler et
78 al., 2021), South America (Alvarez-Garreton et al., 2018; Chagas et al., 2020), and China (Hao et
79 al., 2021) to contribute their regional CAMELS. Another comprehensive regional LSH dataset for
80 North America named the Hydrometeorological Sandbox - École de Technologies Supérieure



Manuscript in submission to *ESSD*

81 (HYSETS) dataset, is also developed with larger sample size (14425 watersheds) and richer data
82 sources compared with the CAMELS (Arsenault et al., 2020).

83 While these datasets provide reliable data sources for regional studies, attempts on building
84 global datasets has become the new norm in the era of big data to boost our analytical and modelling
85 capability for the terrestrial hydrological processes. The HydroATLAS dataset integrates indices of
86 hydrology, physiography, climate, land cover, soil, geology and anthropogenic activity attributes for
87 8.5 million global river reaches (Lehner et al., 2022; Linke et al., 2019). A recent work combined a
88 series of CAMELS dataset with HydroATLAS attributes into a new global community dataset on
89 the cloud named Caravan, with dynamic hydro-climate variables and comprehensive static
90 catchment attributes extracted on 6830 watersheds (Kratzert et al., 2023), which represents by far
91 the most comprehensive synthesis of existing CAMELS. Another global-scale effort, the Global
92 Streamflow Indices and Metadata archive (GSIM), incorporated dynamic streamflow indices and
93 attribute metadata for topography, climate type, land cover, etc., for over 35000 gauges (Do et al.,
94 2018; Gudmundsson et al., 2018), and the streamflow indices are updated to allow for trend analysis
95 (Chen et al., 2023). A recent study filled in the discontinuity and latency of gauge records, and
96 provided streamflow for over 45,000 gauges with improved data quality (Riggs et al., 2023). These
97 global-scale datasets have been widely used in data driven machine learning models (Kratzert et al.,
98 2019a, 2019b; Ren et al., 2020), physical hydrological models (Aerts et al., 2022; Clark et al., 2021),
99 and parameter estimation and regionalization studies (Addor et al., 2018; Fang et al., 2022).

100 Although the flourishing of LSH datasets has promoted comparative hydrological studies
101 (Kovács, 1984) and large-scale hydrological modeling and analyses, several challenges are still
102 standing in the way to realize the full potential of LSH. As briefly outlined in a recent review by
103 Addor et al. (2020), current LSH datasets lack common standards, metadata and uncertainty
104 estimates, and are insufficient in characterising human interventions. More specifically, the
105 following major critical aspects still need attentions from the LSH developers, which we attempt to
106 address with GSHA (Yin et al., 2023). First, the majority of current datasets (especially those at a
107 global scale) incorporate only one data source for each variable, while earth observations, reanalysis,
108 satellite-based estimates are subject to uncertainties (Merchant et al., 2017; Ukhurebor et al., 2020).
109 These uncertainties were rarely represented and may bring difficulties to the regionalization of
110 model parameters (Beck et al., 2016), while also resulting in inconsistent conclusions. Second,
111 anthropogenic activities including land use and land cover (LULC) changes, dam and reservoir
112 building, etc., are critical drivers of shifts in streamflow statistical moments (Niraula et al., 2015).
113 However, historical time series of watershed human modifications have rarely been included in LSH
114 datasets, which is particularly problematic for regions with rapid economic growth. Finally,
115 although the most recent Caravan has provided hydroclimate data for global watersheds, the samples
116 are limited to the existing regional CAMELS which Caravan synthesizes. Therefore, plenty of room
117 is left to increase data sample size and spatial coverage by revisiting the streamflow data acquisition
118 process in a more comprehensive way.

119 To complement existing LSH datasets, we contribute the first version of a synthesis of Global
120 Streamflow characteristics, Hydrometeorology, and catchment Attributes (GSHA v_1.0) for large
121 sample river-centric studies. GSHA features the following characteristics:

- 122 ● Updated physical and anthropogenic descriptors of global rivers, covering streamflow
123 characteristics, hydrometeorological variables, and land use land cover changes for 21568
124 watersheds derived from gauged streamflow records from 13 agencies.



Manuscript in submission to *ESSD*

- 125 ● Streamflow indices for data scarce regions, including those derived from 263 gauges in
 126 China, are included.
 127 ● Extended temporal coverage for as long as 43 years (1979-2021), which varies regionally.
 128 ● Uncertainty estimates for the meteorological variables.
 129 ● Dynamic descriptors for the urban, forest, and cropland fractions, as well as reservoir
 130 storage capacity to improve the representation of human activities in the basin.

131 With the above features, we expect GSHA to support hydrological model parameter estimation
 132 and data driven analysis of global streamflow as one of the most comprehensive LSH datasets
 133 regarding sample size, variable dynamics, and uncertainty estimates. **Table 1** summarizes the
 134 differences between GSHA and other prominent LSH datasets. Our paper is organized as follows.
 135 Section 2 expands on **Table 1** and provide more details of the data included for GSHA. Section 3
 136 introduces the data sources and methodologies involved in creating GSHA. Section 4 highlights the
 137 key features of GSHA by conducting some analyses, followed by conclusions reached in Section 5.
 138

139 **Table 1 Comparison of GSHA with other LSH datasets.** Note that we only include the CONUS
 140 CAMELS dataset to represent regional LSH datasets for this comparison, as other regional CAMELS
 141 share large similarity with CONUS CAMELS.

Factors	CAMELS (eg. US)	HydroATLAS	Caravan	GSIM	GSHA
Spatial extent	Regional	Global	Global	Global	Global
Sample size	671	8.5 million	6830	35002	21568
Time span	1980–2015	Static	1981–2020	1806-2016	1979-2021
Streamflow dynamics	Yes	No	Yes	Yes (statistical indices)	Yes (statistical indices)
Meteorological time series	Yes	No	Yes	No	Yes
Multi data sources for meteorological variables	Yes	No	No	No	Yes (with uncertainty estimates)
Water storage dynamics	No	No	Only soil water dynamics	No	Yes
Land cover dynamics	No	No	No	No	Yes
Reservoir dynamics	No	No	No	No	Yes
Static attributes	Yes	Yes	Yes (from HydroATLAS)	Yes	Yes (from HydroATLAS)



Manuscript in submission to *ESSD*

142 2 Dataset content of GSHA v1

143 In this section, the data fields, variables, and attributes included in GSHA are described in more
144 details and summarized in **Table 2**. For the instructions of the data format, we provided a user
145 manual along with the dataset (see `readme.docx`). GSHA includes yearly streamflow characteristics
146 derived from daily discharge observations, meteorological variables (including precipitation, 2-m
147 air temperature, long- and shortwave radiation, wind speed, actual and potential evapotranspiration
148 (AET and PET)), daily or weekly water storage terms (4 layers of soil moisture, groundwater, and
149 snow depth water equivalence), daily vegetation index (leaf area index (LAI)), yearly LULC
150 characteristics (urban, cropland, and forest fraction), and yearly reservoir information (degree of
151 regulation (DOR) and reservoir capacity). For each meteorological variable, multiple independent
152 data sources are incorporated to provide uncertainty estimates. Static attributes like land
153 physiography, soils, and geology are not additionally extracted, as similar efforts have been made
154 by other researchers, so we directly matched our gauge locations to the HydroATLAS dataset
155 (Lehner et al., 2022; Linke et al., 2019) by providing the river ID match table. Users can link the
156 two to obtain these attributes.

157 **Watershed polygons:** GSHA includes 21568 watershed polygons delineated from the global
158 gauges, which is stored as Esri Shapefile format. The ID and agency of each watershed is the same
159 as the corresponding gauge ID, and the gauge latitude/longitude are in decimal degree. The area
160 denotes the upstream drainage basin area of the gauge. Some of the IDs contain characters (such as
161 '.', '-', etc.) inconsistent with the majority of IDs. For the convenience of the users, we unified these
162 as underscores and stored the new file names as 'filename'.

163 **Streamflow indices:** GSHA publishes annual streamflow indices derived from daily
164 streamflow data, including different percentiles, and mean/median/minimum/maximum. The
165 frequency and durations of extremely high and low streamflow events are also provided, along with
166 numbers of zero observations and valid samples to allow extreme streamflow analysis and flexible
167 data screening by the users. The indices are stored as comma separated values (CSV) files by year,
168 with each watershed corresponding to one file. A complementary R package can be used to
169 automatically download many of the gauge datasets is available at [https://github.com/Ryan-](https://github.com/Ryan-Riggs/RivRetrieve)
170 [Riggs/RivRetrieve](https://github.com/Ryan-Riggs/RivRetrieve)(Riggs et al., 2023).

171 **Meteorological variables:** The meteorological variables selected are the most influential
172 drivers for streamflow, which includes precipitation, 2-m temperature, ET, radiation and wind speed.
173 In main-stream land surface models, ET is a diagnostic variable derived from meteorological inputs
174 and is not considered as meteorological forcing. However, as many hydrological models also use
175 potential ET as an input variable, and model calibration sometimes involves actual ET (Immerzeel
176 & Droogers, 2008), we include the two variables and place it into the meteorological variable
177 category. For each variable, more than one data sources are used to allow for uncertainty analysis,
178 which will be provided on a yearly basis in an independent file.

179 **Natural water storage terms and land use/land cover change:** These include soil moisture,
180 snow water equivalent, and ground water percentages. We also include yearly land cover dynamics



Manuscript in submission to *ESSD*

181 (i.e., urban, forest, and cropland fraction changes), as well as dynamically changing reservoir
 182 capacity and degree of regulation (DOR) percentage. Leaf area index (LAI) is also included to
 183 reflect the seasonal changes in vegetation canopy that is also key to the streamflow processes.

184 **Static attributes:** GSHA does not extract updated static attributes because HydroATLAS
 185 already make substantial efforts in this regard. Instead, the listed categories are those mostly related
 186 to streamflow prediction from HydroATLAS selected to be included in GSHA files, and we direct
 187 the readers to the ID match table to access the entire 281 static attributes offered by HydroATLAS
 188 (Lehner et al., 2022; Linke et al., 2019). Our user manual, available at the dataset download site,
 189 also provides more information on it.

190

191 **Table 2 Fields provided with GSHA.**

Category	Field	Description	Unit
Watershed	Sttn_Nm	The ID of the watershed.	NaN
Polygons	Latitude	Latitude of the gauge.	Degree
	Longitude	Longitude of the gauge.	Degree
	Shedarea	The area of delineated watershed.	Km ²
	Agency	The agency the gauge belongs to.	NaN
	filename	The name of the corresponding Shapefile in the dataset.	NaN

Category	Indices	Description	Unit/Format
Streamflow indices	percentiles	Annual 1, 10, 25, 75, 90, 99 percentiles of daily streamflow.	m ³ /s
	mean	Annual mean of daily streamflow.	m ³ /s
	median	Annual median of daily streamflow.	m ³ /s
	annual maximum flood (AMF)	Annual maximum of daily streamflow.	m ³ /s
	AMF occurrence date	The date of AMF occurrence.	Year/month/day
	frequency of high-flow events	Number of days in a year with streamflow \geq 90 percentile flow.	Days/year
	average duration of high-flow events	Average number of consecutive days \geq 90 percentile flow.	Days
	frequency of low-flow events	Number of days in a year with streamflow \leq 10 percentile flow.	Days/year
	average duration of low-flow events	Average number of consecutive days \leq 10 percentile flow.	Days
	Q=0 days	Number of days with runoff=0.	Days
valid observation days	Number of days with no missing data. (Valid observations refer to non-null measurements.)	Days	

Category	Variable	Data source name	Unit
Meteorological	Precipitation	MSWEP	mm
Variables		EM-Earth	mm
	2 m temperature	ERA5	K



Manuscript in submission to *ESSD*

		MERRA-2	K
		EUSTACE	K
Actual		REA	mm
evapotranspiration		GLEAM	mm
Potential		GLEAM	mm
evapotranspiration		hPET	mm
Radiation (longwave)		ERA5 land surface net thermal radiation	W/m ²
		MERRA-2 surface net downward longwave flux	W/m ²
Radiation (shortwave)		ERA5 land surface net solar radiation	W/m ²
		MERRA-2 surface net downward shortwave flux	W/m ²
10 m wind speed (u component)		ERA5 land u-component of wind	m/s
		MERRA-2 10 metre eastward wind	m/s
10 m wind speed (v component)		ERA5 land v-component of wind	m/s
		MERRA-2 10 metre northward wind	m/s
10 m wind speed (actual)		ERA5 land u- and v-components of wind	m/s
		MERRA-2 10 metre northward and eastward wind	m/s
Category	Variable	Data source name	Unit
Water storage terms	Soil moisture layer 1	ERA5 land soil water layer 1 (0-7 cm, 0cm refers to the surface)	m ³ /m ³
	Soil moisture layer 2	ERA5 land soil water layer 2 (7-28 cm)	m ³ /m ³
	Soil moisture layer 3	ERA5 land soil water layer 3 (28-100 cm)	m ³ /m ³
	Soil moisture layer 4	ERA5 land soil water layer 4 (100-289 cm)	m ³ /m ³
	Snow water equivalent	ERA5 land snow depth water equivalent	m of water equivalent
	Ground water	GRACE-FO data assimilation	%
Category	Variable	Data source name	Unit
Land use and land cover	Urban fraction	GAUD	%
	Forest fraction	MCD12Q1	%
	Cropland fraction	MCD12Q1	%
	Reservoir capacity	GeoDAR	Million m ³
	DOR	GeoDAR	%
LAI	CDR LAI	NaN	
Category	Attribute	Column name	Unit
Static-	Elevation	ele_mt_uav	m. a.s.l.
Physiography	Terrain slope	slp_dg_uav	degrees (x10)
	Stream gradient	sgr_dk_rav	decimetres per



Manuscript in submission to *ESSD*

				km
Static-	Inundation Extent		inu_pc_ult	%
Hydrology	Groundwater	Table	gwt_cm_cav	cm
	Depth			
Static-	Land Cover Classes		glc_cl_cmj	NaN
Landcover	Potential	Natural	pnv_cl_cmj	NaN
	Vegetation Classes			
	Wetland Extent		wet_pc_u01-u09	%
	Glacier Extent		gla_pc_use	%
	Permafrost Extent		prm_pc_use	%
Static-Soil &	Clay Fraction in Soil		cly_pc_uav	%
geology	Silt Fraction in Soil		slt_pc_uav	%
	Sand Fraction in Soil		snd_pc_uav	%
	Lithological Classes		lit_cl_cmj	NaN
	Soil Erosion		ero_kh_uav	kg/hectare per year

192 3 Data sources and methodology

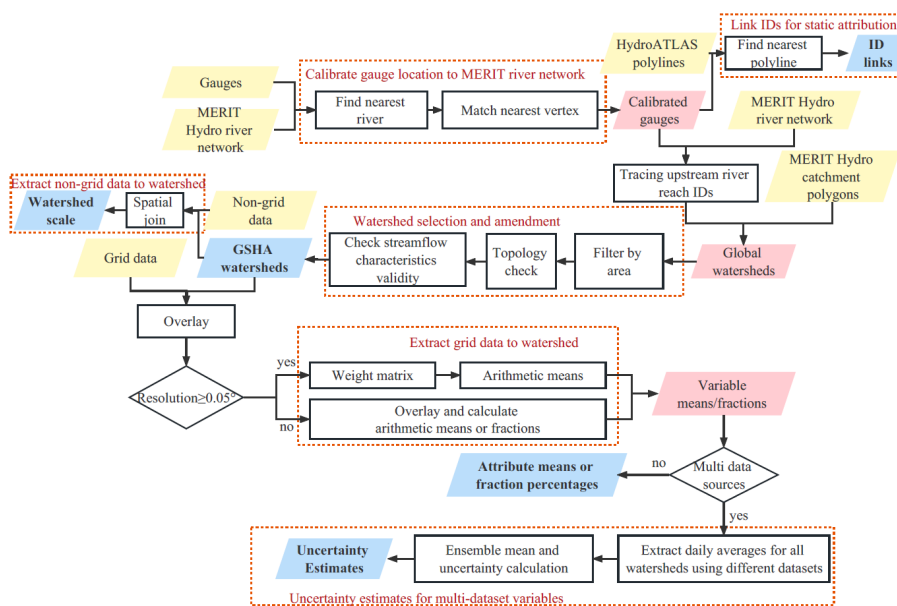
193 3.1 Technical workflow in creating GSHA

194 The creation of GSHA starts from revisiting the data compilation process for the stream
 195 gauging observations from 13 international agencies. The general workflow of GSHA data
 196 production processes is illustrated in **Figure 1**, which consists of watershed delineation, variable
 197 extraction from both grid and non-grid data sources, and uncertainty analysis.

198 First, we delineated the upstream watersheds using gauge locations. Calibration of gauge
 199 longitudes and latitudes were conducted to match the gauges with the MERIT river network exactly.
 200 The delineated watersheds were selected and manually checked using standards of area, topology
 201 correctness, and observation data lengths. The selected watersheds went on to be overlaid with
 202 grid and non-grid variable data sources for to obtain GSHA variables.



Manuscript in submission to *ESSD*



203
 204
 205
 206
 207

Figure 1 General workflow of GSHA. The yellow parallelograms are the input datasets, the blue ones are the final outputs of GSHA dataset, and the pink ones are the results in the process. The black quadrilaterals represent the extraction and calculation processes, and the red dotted rectangles illustrate different modules of the extraction process.

208 3.2 Gauge-based streamflow indices

209 As shown in **Table 3**, in total 36497 gauges were initially scraped from web and from the
 210 Chinese National Real-time Rain and Water Situation Database including the real-time water level
 211 and streamflow data of the hydrological stations. For gauges located within ~100-m of each other,
 212 those with shorter record lengths are removed, assuming that they are redundant with one another.
 213 The gauge measurements were converted to a consistent unit (m^3/s) and were then manually
 214 compared with GRDC measurements to ensure accurate unit conversion (Riggs et al., 2023). Gauge
 215 databases compiled in this study are available through a variety of web interfaces, except for the
 216 CHP data which is provided by the authors of the dataset (Henck et al 2010, Schmidt et al 2011),
 217 and processed into annual scale data that meets the requirements of the synthesis dataset.



Manuscript in submission to *ESSD*

Table 3 Gauge data sources used in this analysis. N1 and N2 refers to numbers of gauges with observations after 1979 and used in GSHA. The starting and ending years (Y1 and Y2) of GSHA gauges for each agency are listed.

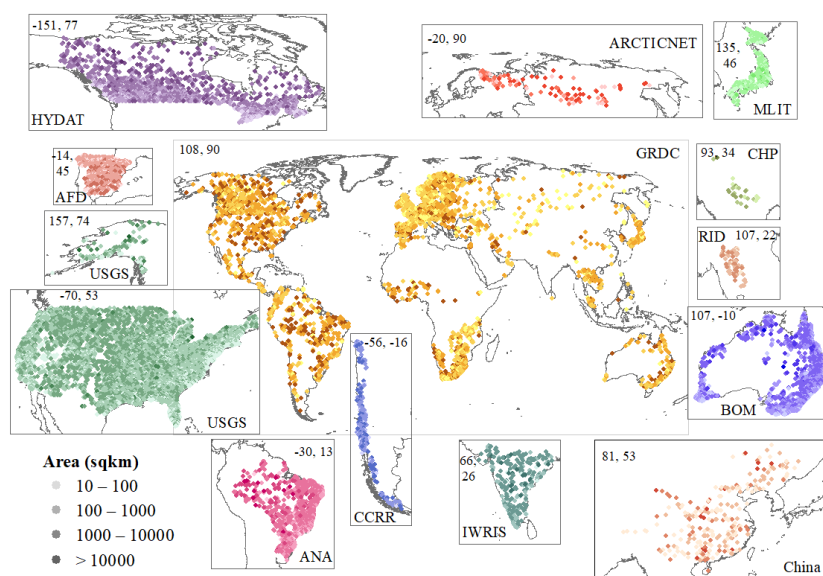
Source	N1	N2	Y1	Y2	URL/Provider
ArcticNET 2022	116	106	1979	2003	www.r-arcnetcnet.sr.unh.edu/v4.0/AllData/index.html
Australian Bureau of Meteorology 2022 (BOM)	4017	2340	1979	2021	www.bom.gov.au/waterdata/
Brazil National Water Agency 2022 (ANA)	1343	1172	1979	2021	www.snih.gov.br/hidroweb/serieshistoricas
Canada National Water Data Archive 2022 (HYDAT)	3771	2222	1979	2021	www.canada.ca/en/environment-climate-change/services/water-overview/quantity/monitoring/survey/data-products-services/national-overview/
Chile Center for Climate and Resilience Research 2022(CCRR)	481	392	1979	2020	https://explorador.c2.cl/
Chinese Hydrology Project (CHP)	112	26	1979	1987	(Henek et al 2010, Schmidt et al 2011)
The Global Runoff Data Centre 2022 (GRDC)	6345	4004	1979	2021	(https://portal.grdc.bafg.de/applications/public.html?publicuser=PublicUser)
India Water Resources Information System 2022 (WRIS)	547	261	1979	2020	https://indiawrts.gov.in/wris/#/RiverMonitoring
Japanese Water Information System 2022 (MLIT)	1023	751	1979	2019	www.lriver.go.jp/
Spain Anuario de Aforos, 2022 (AFD)	1138	889	1979	2018	http://datos.gob.es/catalogo/00125801-anuario-de-aforos/resource/4836b826-e7fd-4a41-950c-89b4eca0279
Thailand Royal Irrigation Department 2022 (RID)	126	73	1980	1999	http://hydro.iis.u-tokyo.ac.jp/GAME-T/GAIN-T/routine/rid-river/disc_d.html
U.S. Geological Survey 2022 (USGS)	16951	9069	1979	2021	https://waterdata.usgs.gov/nwis/rt
Chinese National Real-time Rain and Water Situation Database	527	263	2000	2019	http://xxfb.mwr.cn/sq_zdysq.html



Manuscript in submission to *ESSD*

218 3.2 Watershed delineation

219 The watershed delineation process is built upon a vector-based global river network dataset
220 (Lin et al., 2021), which is delineated from the 90-m Multi-Error-Removed Improved Terrain
221 (MERIT) digital elevation model (DEM) (Yamazaki et al., 2017) and the flow direction and flow
222 accumulation rasters (Yamazaki et al., 2019). The locations of the gauges may contain locational
223 errors and direct delineation will result into erroneous watershed boundaries; therefore, gauge
224 location correction was conducted by relocating the gauges to the nearest MERIT-based river reach
225 vertices. The adjusted gauge points are used as the watershed outlets, where the contributing areas
226 are extracted by dissolving all upstream catchments based on the topology provided by MERIT
227 Basins (Lin et al., 2019). Since the area threshold of MERIT Basins is 25 km², we do not include
228 watersheds smaller than this threshold. Considering the spatial heterogeneity of very large basins,
229 we excluded watersheds $\geq 50,000$ km² from the dataset. To ensure GSHA to support studies with
230 sufficiently long records, only watersheds with >5 years of observations since 1979 were selected.
231 For gauges sharing the same watershed, the one with better data quality (i.e., longer measurement
232 records and more valid observation days) is used. If the two gauges share the same quality, we only
233 included the furthest downstream gauge. Eventually, the selection processes resulted in 21568 valid
234 watersheds out of 35970 gauges initially scraped from web plus 527 gauges from the Chinese
235 National Real-time Rain and Water Situation Database (**Figure 2**).



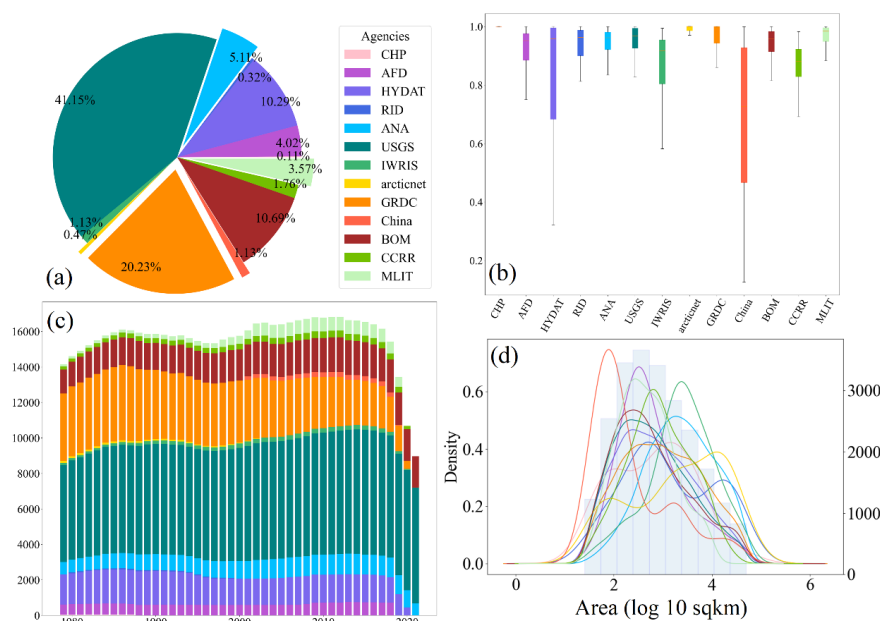
236
237 **Figure 2 Spatial distribution of the GSHA gauges (n=21568).** Watershed areas are represented by the
238 tint of colours. Gauges of different agencies are represented with separate colours and are plotted in
239 individual frames (except for USGS gauges in two frames to incorporate Alaska). The agency names and
240 the upper-left coordinates (longitude, latitude) of each frame are also shown in the figure.

241
242 The GSHA watersheds are unevenly distributed across the globe, more than half of which are
243 located in North America (USGS, HYDAT and a large proportion of GRDC gauges, **Figure 3a**).



Manuscript in submission to *ESSD*

244 Europe, Australia and South America also have a relatively good coverage, while Asia and Africa
 245 show the lowest gauge densities. The majority of the gauged watersheds are of medium sizes ranging
 246 from 250 to 2500 km², although for some agencies it does not show the same distribution (**Figure**
 247 **3d**). For instance, ANA (South America), IWRIS (India) and arcticnet (Northern Eurasia)
 248 watersheds are generally larger, while the Chinese National Real-time Rain and Water Situation
 249 Database provides more gauges with smaller drainage areas. Due to the maintenance difficulties,
 250 the numbers of functioning gauges are declining for agencies like GRDC, but the lack of data in
 251 recent years (**Figure 3c**) are mainly due to latency issues. USGS, BOM and ANA provide a stable
 252 number of observations for the 1980-2021 period (**Figure 3c**) with high proportions of valid
 253 observations each year (**Figure 3b**), while observational periods from arcticnet and China contain
 254 relatively fewer valid samples (**Figure 3b**) and shorter time spans (**Figure 3c**).



255
 256 **Figure 3 Summary statistics of the GSHA gauges.** This includes (a) proportions of gauges from
 257 different agencies, (b) box plots for proportions of valid observations for each agency, (c) proportion of
 258 valid observation for each year by agency and (d) distributions of watershed areas for each agency
 259 (kernel density estimation lines, left y-axis) and all gauges (blue histogram, right y-axis). The colour
 260 legend in subplot (a) applies to all four subplots. In subfigure (a) the 0.11% label corresponds to CHP,
 261 and the legend goes counter clockwise in the pie chart. In subfigure (c), CHP bars are at the bottom of the
 262 plot, and the legend goes from bottom to the top of the bars.

263 3.3 Meteorological variables, water storage terms, and land surface characteristics

264 After watershed delineation, publicly available grid or non-grid data were obtained and
 265 overlaid to derive the meteorological, water storage terms, and land surface characteristics. The data
 266 sources used for GSHA are listed in **Table 4**. We prioritized the use of multi-source fusion datasets



Manuscript in submission to *ESSD*

267 with relatively high quality surveyed from literature when creating GSHA.

268 3.3.1 Meteorology datasets

269 For precipitation, the Multi-Source Weighted-Ensemble Precipitation (MSWEP) that have
270 merged gauge (CPC Unified), grid (GPCC), satellite (CMORPH, GSMaP-MVK, and TMPA
271 3B42RT), and reanalysis data (ERA-Interim and JRA-55) with sample density and comparative
272 performance considered (Beck et al., 2017; Beck et al., 2019) is used. Another recent precipitation
273 dataset is the Ensemble Meteorological Dataset for Planet Earth (EM-Earth) deterministic estimates,
274 which merges a station-based Serially Complete Earth (SC-Earth) removing the temporal
275 discontinuities in raw station observations and ERA5 estimates (Tang et al., 2022). The EUSTACE
276 global land station daily air temperature dataset (EUSTACE) statistically merged station and
277 satellite observations to obtain global daily near-surface air temperature (Brugnara et al., 2019), and
278 is used as a source for 2-m temperature in GSHA. Other datasets used for 2-m temperature extraction
279 are the reanalysis datasets Modern-Era Retrospective analysis for Research and Applications
280 Version 2 (MERRA-2) (Gelaro et al., 2017) by NASA's Global Modelling and Assimilation Office
281 (GMAO) and the land components of the European Centre for Medium-Range Weather Forecasts
282 (ECMWF) fifth generation of European Reanalysis (ERA5) dataset (Muñoz-Sabater et al., 2021).
283 These reanalysis datasets are also used in extracting long- and shortwave radiation, as well as u- and
284 v-components of wind. MERRA-2 uses the Goddard Earth Observing System (GEOS) model and
285 analysis scheme, and assimilated the latest observations. The land surface model used in ERA5
286 reanalysis is the Carbon Hydrology-Tiled ECMWF Scheme for Surface Exchanges over Land
287 (CHTESSEL) driven by the downscaled meteorological forcing from the ERA5 climate reanalysis
288 (Hersbach et al., 2020). For AET, the dataset merging ERA5, Global Land Data Assimilation System
289 Version 2 (GLDAS2), and MERRA-2 using the reliability ensemble averaging (REA) method is
290 used (Lu et al., 2021), together with the product of Global Land Evaporation Amsterdam Model
291 (GLEAM) based on satellite observations of surface net radiation and near-surface air temperature
292 (Martens et al., 2017). For PET, GLEAM is also incorporated. Another PET dataset for GSHA
293 is the hourly PET at 0.1° resolution for the global land surface from 1981-present (hPET) dataset
294 calculated from ERA5-land wind speed, air and dew point temperature, net radiation components
295 and surface air pressure (Singer et al., 2021).

296 3.3.2 Water storage term datasets

297 ERA5-land data was also applied in extracting soil moisture for 4 soil layers, as well as snow
298 water equivalence. For groundwater, an assimilation dataset from the NASA's Gravity Recovery
299 and Climate Experiment (GRACE) and its follow-on mission (GRACE-FO) is used (Li et al., 2019).
300 The dataset merged water storage derived from GRACE satellite products into ECMWF Integrated
301 Forecasting System meteorological data-forced NASA's Catchment land surface model (CLSM).
302 The data is represented as groundwater drought indicator (GWI), which is the percentage of
303 groundwater storage estimates from the GRACE data assimilation relative to the climatology
304 (representing historical conditions), at weekly time scales from 2003-2021.



Manuscript in submission to *ESSD*

305 3.3.3 Land surface characteristic datasets

306 Global urban development for 1985-2015 is represented as the urban fraction in each watershed
307 using the global annual urban dynamics (GAUD) at 30-m resolution. The dataset was derived from
308 Landsat surface reflectance based on the Normalized Urban Areas Composite Index (NUACI) (Liu
309 et al., 2020). For forest and cropland fractions, the Terra and Aqua combined Moderate Resolution
310 Imaging Spectroradiometer (MODIS) Land Cover Type (MCD12Q1) land cover dataset, was used
311 (Friedl et al., 2010). It covers 2001-2020 with a resolution of 500 m, and the categories used for
312 GSHA is the International Geosphere–Biosphere Programme classification (IGBP) forests and
313 croplands. Another land cover is vegetation, which is represented by LAI obtained from the National
314 Oceanic and Atmospheric Administration (NOAA) Climate Data Record (CDR) of Advanced Very
315 High-Resolution Radiometer (AVHRR) product, which relies on artificial neural networks and
316 AVH09C1 surface reflectance product (Claverie et al., 2016).

317 3.3.4 Dams and reservoirs

318 The newly published Georeferenced global Dams And Reservoirs (GeoDAR) dataset that
319 documents the dam and reservoir construction years is used for building the temporally varying
320 watershed reservoir capacity and DOR. GeoDAR georeferenced the International Commission on
321 Large Dams (ICOLD) World Register of Dams (WRD), and geo-matched multi-source regional
322 registers and geocoding descriptive attributes through the Google Maps API (Wang et al., 2022).
323 The reservoir capacities are used together with the mean annual streamflow to obtain the DOR based
324 on the equation $dor = SC/Q_{mean}$, where SC refers to reservoir storage capacity and Q_{mean} is
325 the mean annual streamflow in the corresponding year.

326 3.3.5 Static variables

327 We matched GSHA river IDs and HydroATLAS river reach IDs to link the static attributes.
328 HydroATLAS includes 56 variables for hydrology, physiography, climate, land cover & use, soils
329 & geology, and anthropogenic influences for over 8.5 million river reaches globally.

330

331 **Table 4 Data sources used for the GSHA variables.**

Category	Dataset	Resolution	Interval	Reference
Meteorology	MSWEP	0.25°	Daily	(Beck et al., 2017; Beck et al., 2019)
	EM-Earth	0.1°	Daily	(Tang et al., 2022)
	ERA5-land	0.1°	Hourly	(Muñoz-Sabater, 2019)
	MERRA-2	0.5°* 0.625°	Hourly	(GMAO, 2015)
	EUSTACE	0.25°	Daily	(Brugnara et al., 2019)
	REA	0.25°	Daily	(Lu et al., 2021)
	GLEAM	0.25°	Daily	(Martens et al., 2017; Miralles et al., 2011)



Manuscript in submission to *ESSD*

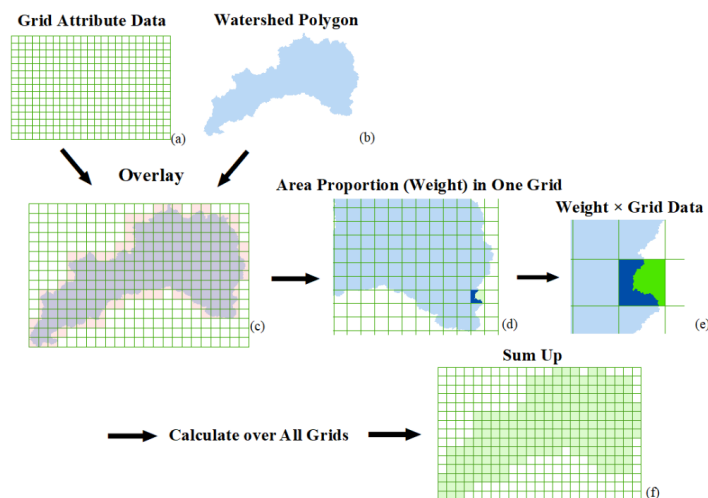
	hPET	0.1°	Daily	(Singer et al., 2021)
Water storage terms	ERA5-land	0.1°	Hourly	(Muñoz-Sabater, 2019)
	GRACE-FO data assimilation	0.25°	Weekly	(Li et al., 2019; Zaitchik et al., 2008)
Land surface	GAUD	30 m	Yearly	(Huang, 2020)
	MCD12Q1	500 m	Yearly	(Friedl et al., 2019)
	CDR Leaf Area Index	0.05°	Daily	(Vermote et al., 2019)
Dam and reservoir	GeoDAR	NaN (polygon)	Yearly	(Wang et al., 2022)
Static Attributes	HydroATLAS	NaN (line)	NaN (static)	(Lehner et al., 2022; Linke et al., 2019)

332 3.4 Variable extraction methods

333 For grid data with relatively coarse spatial resolutions ($\geq 0.05^\circ$), we used an area-weighted
 334 approach to calculate the weighted average of the variable, while for high-resolution grid data, we
 335 extract the arithmetic mean directly. **Figure 4** shows the area-weighted average approach we used
 336 for grid data with spatial resolution $\geq 0.05^\circ$ to reduce the influence of watershed area on data
 337 uncertainty (Tang et al., 2022). The grid data (**4a**) and the quality-controlled watersheds (**4b**) are
 338 overlaid and all grids intersecting with the watershed are obtained (**4c**). For each intersected grid,
 339 the proportion of the polygon in the grid is calculated as the weight (dark blue, **4d**); the product of
 340 the weight and the corresponding grid value is calculated over all intersected grids (**4e**) and are
 341 summed up as the weighted average (**4f**). For wind, the u- and v-wind components are first used to
 342 calculate wind speed, then the basin average is calculated with the weighted average approach. For
 343 grid data with a spatial resolution of $< 0.05^\circ$, the area-weighted approach is not adopted as it offers
 344 limited gains while becoming computationally too expensive. For reservoirs, they are spatially
 345 joined to GSHA watershed polygons, and all the intersected reservoirs were used to calculate the
 346 total reservoir storage capacity and degree of regulation.



Manuscript in submission to *ESSD*



347

348

349

350

Figure 4 Determination of the area weights in extracting gridded data to GSHA watershed polygons. This weighted approach is applied to data at a resolution of $\geq 0.05^\circ$ but not for data at a finer spatial resolution due to computational costs.

351

3.5 Uncertainty estimates

352

For meteorological variables, uncertainty estimates are calculated using Eq. (1):

353

$$\text{uncertainty} = \frac{X_{max} - X_{min}}{\bar{X}} * 100\%, \quad (1)$$

354

355

356

357

358

where X_{max} and X_{min} are the maximum and minimum among the extracted values from the independent data sources. \bar{X} is the mean of values from all datasets. The uncertainty ranges from 0 to 200%. Note that for some of our data sources such as EC-Earth, uncertainty estimates are intrinsically provided, but here we do not include it into GSHA as we consider the uncertainty estimates made from independent data sources more internally consistent among different variables.

359

3.6 Validation

360

361

362

363

364

365

366

367

368

369

Postprocessing of the extracted variables include the unification of units and manual quality checks. For streamflow characteristics, we validated three of our indices against GSIM for its global coverage, including the mean annual streamflow, 10th and 90th percentiles. The spatial joint between GSHA and GSIM gauges in a 10 km buffer zone was performed, and only the GSIM gauge with a minimum distance and watershed area difference $\leq 5\%$ to a GSHA gauge is considered. Pairs with 0 measurements were excluded and 9835 pairs were involved eventually. We plotted the scatter plot of GSHA-GSIM mean flow, 10th and 90th percentiles, and compared the fitting line to the 1:1 line, with correlation coefficients calculated (see Section 4.1).

We also validated precipitation, potential ET and 2 m air temperature with the regional CAMELS-US dataset. We compare the Daymet meteorological variables of CAMELS and the mean



Manuscript in submission to *ESSD*

370 of GSHA variables for the validation. Since we include ERA5 data for most of our variables directly
371 or indirectly as the data source, while Caravan consistently used ERA5, we did not use Caravan for
372 the global validation as it is not considered as fully independent from GSHA. The spatial match is
373 the same we did for GSIM which resulted in 906 pairs. This number is larger than the total CAMELS
374 gauge numbers as some gauges might be repeatedly paired due to location bias of the USGS gauges
375 and MERIT river networks, as well as the adjacency between gauges of different agencies. Similarly,
376 scatter plots and correlation coefficients are provided for assessment.

377 3.7 Watershed classification and change detection

378 We classified the watersheds as natural and human managed to analyse the influence of human
379 water management. A watershed is classified as a natural watershed if it satisfies the following: (1)
380 DOR is smaller than 10%; (2) the urban extent is less than 5%; and (3) the sum of urban and cropland
381 fractions is smaller than 10% (L. Yang et al., 2021; Zhang et al., 2023). The classification was
382 performed for 2001-2015, and the changing patterns of the watersheds are divided into six categories:
383 (1) natural (N) when the watershed remained natural for all 15 years; (2) human managed (H) when
384 the watershed remained human managed for all 15 years; (3) natural to human managed (NH) when
385 the watershed was first natural in 2001, but changed to and maintained human managed later; and
386 (4) human managed to natural (HN) when the watershed was first human managed in 2001, but
387 changed to and maintained natural later.

388 4 Results

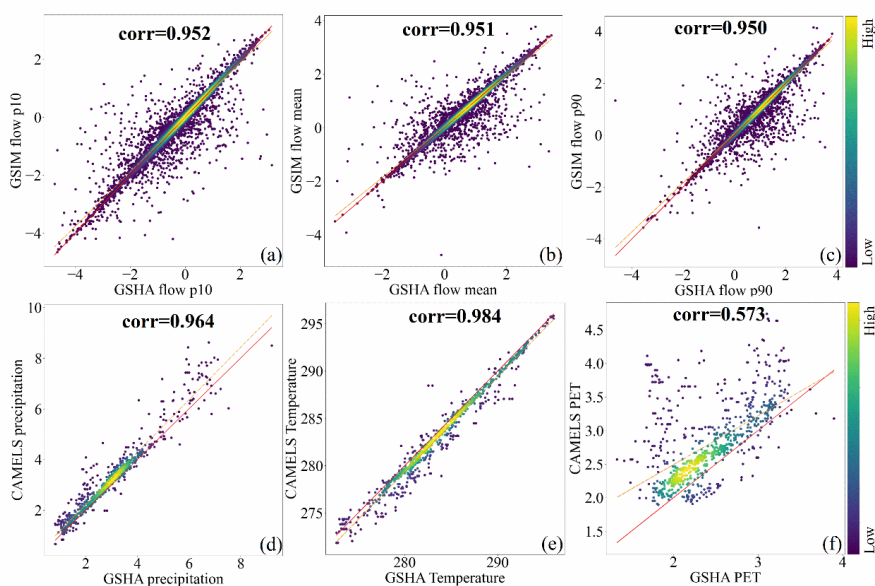
389 As previous studies have already revealed the spatial patterns of the LSH hydrometeorological
390 variables both locally and globally, therefore here we put the spatial patterns of GSHA
391 meteorological variables and streamflow indices in **Appendix A**, while we focus on using the
392 Results section to reveal the uniqueness of GSHA, which includes a technical validation of GSHA,
393 uncertainty analysis, and the temporal change of watershed human management levels.

394 4.1 Technical validation

395 **Figure 5** illustrates the validation results of GSHA. **Figures 5a–5c** show streamflow indices
396 as validated against GSIM globally, and **Figures 5d–5f** show meteorological variable as validated
397 against Daymet from CONUS CAMELS. For streamflow indices, precipitation, and temperature,
398 the correlation coefficients exceed 0.95 (significance $p < 0.01$), and the fitting lines are close to 1:1
399 line, indicating high consistencies between GSHA and the reference datasets. For PET, however, the
400 coefficient is low, at only 0.573 (significance $p < 0.05$), and the CAMELS PET is generally higher
401 than GSHA ensemble, which is possibly ascribed to the high uncertainty among PET datasets that
402 is yet to be fully resolved (Singer et al., 2021) (see **Appendix B**).



Manuscript in submission to *ESSD*



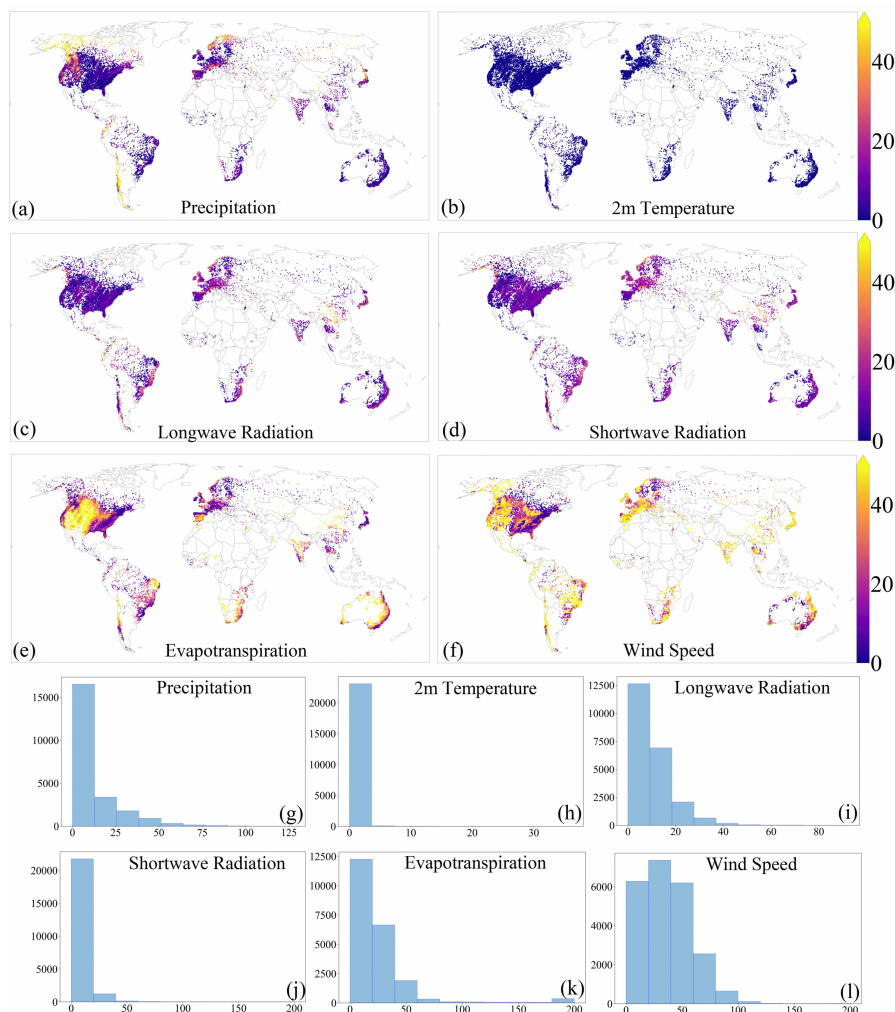
403
404 **Figure 5** Validation of GSHA with GSIM streamflow characteristics ((a), (b) and (c)), and
405 **CAMELS meteorological variables ((d), (e) and (f))**. ‘Corr’ in the subfigure is the Pearson correlation
406 coefficient. The red line is the 1:1 line, while the orange dotted line is the fitting line of the scatter points.
407 The colour bar represents density of the sample points.

408 4.2 Uncertainty patterns for the GSHA meteorological variables

409 **Figure 6** shows the distributions of the uncertainties for different variables, and the colour bars
410 are unified to allow for comparisons between different variables.



Manuscript in submission to *ESSD*



411

412 **Figure 6 Global patterns of the uncertainty for the GSHA meteorological variables (in percentage).**

413 This includes the uncertainty (a) for precipitation (mm/day), (b) 2-m temperature (K), (c) longwave
414 radiation (W/m^2), (d) shortwave radiation (W/m^2), (e) evapotranspiration (mm/day), and (f) wind speed
415 (m/s), and (g) the uncertainty histogram for precipitation, (h) 2-m temperature, (i) longwave radiation,
416 (j) shortwave radiation, (k) evapotranspiration, and (l) wind speed.

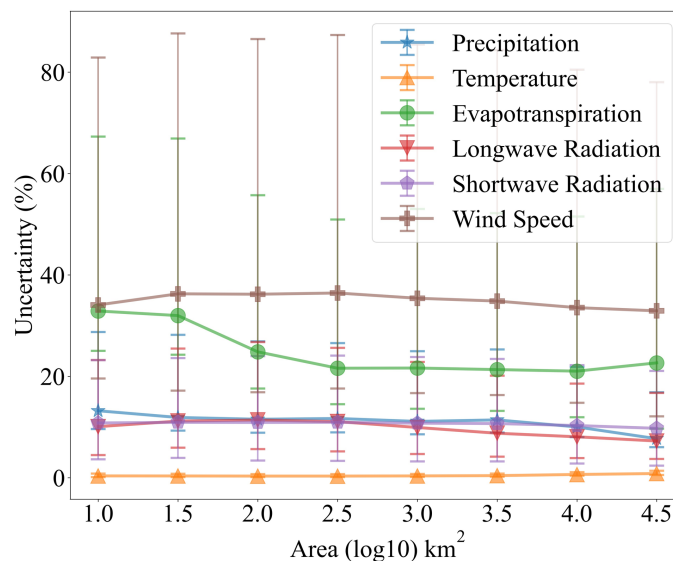
417

418 Generally, among all variables, air temperature (**Figures 6b & 6h**) shows the minimum
419 uncertainty (<5%), while the uncertainty for wind speed (**Figure 6f**) is the highest among all
420 variables. Uncertainties for other variables show strong spatial variability. For example,
421 uncertainties for precipitation are high in high-latitude or mountainous areas like Rocky Mountains,
422 northern Europe, the Alps and the Andes areas (**Figure 6a**). This is reasonable because limited
423 accessibility to in-situ observations and the misestimation of snow (Schreiner-McGraw & Ajami,
424 2020) can contribute to the precipitation estimation errors, while the data sources show relatively



Manuscript in submission to *ESSD*

425 high consistency (*uncertainty* $\leq 25\%$) in other parts of the world (**Figure 6g**). For radiation, as
426 solar/shortwave radiation is largely affected by sky conditions, thus uncertainties are high in regions
427 with less clear sky, including south-west China and its surrounding areas, high latitude regions of
428 the northern hemisphere, and Europe (Brun et al., 2022). These places are also subject to high
429 thermal/longwave radiation uncertainties for similar reasons (**Figure 6c**). Land cover including
430 vegetation and artificial surface, is another factor influencing surface net radiation through albedo
431 effect (Hu et al., 2017), thus for heavily vegetated and urbanized areas, such as the Amazon region
432 and east coastal Australia, uncertainties for both longwave and shortwave fluxes are also relatively
433 high. Nevertheless, **Figures 6i & 6j** demonstrate that for the majority of watersheds, radiation
434 uncertainties are $< 25\%$, indicating that the radiation data sources are generally consistent to each
435 other. ET uncertainties are generally larger than the above variables (**Figures 6e & 6k**), and are
436 particularly prominent in dry areas of the globe, e.g., central North America, northern Andes, central
437 Asia, and Australia's grasslands and deserts. It is also prominent in agriculture intensive regions like
438 India and northern part of China (Sörensson & Ruscica, 2018), where the agricultural irrigation may
439 be the contributing factor to the ET uncertainty. The spatial distributions of wind speed do not seem
440 to show clear regional patterns (**Figure 6f**), and uncertainty values of wind speed are generally
441 larger over majority of watersheds (**Figure 6l**). Nevertheless, the uncertainties are low in the
442 Appalachia and northern Europe, and are high in most parts of Brazil, the Andes, Africa, eastern
443 and southern parts of Asia, as well as Australia (**Figure 6f**). As we already selected relatively high-
444 quality datasets for the variables, these areas might be calling for more attention by the LSH
445 developers, while providing possible explanations for the inconsistencies in interpreting results or
446 understanding the challenges in estimating model parameters by the LSH users.



447
448 **Figure 7 Relationship between variable uncertainties and watershed areas.** The markers indicate
449 mean values of the variable uncertainties in watersheds smaller than the corresponding x-axis value. The
450 error bars represent the range between 25 and 75 percentiles of the uncertainty values.
451



Manuscript in submission to *ESSD*

452 Apart from the spatial characteristics above, we also investigate the emergent patterns of the
453 uncertainties. Existing studies indicate small basins can show larger uncertainties due to coarse resolution
454 data inputs (Kauffeldt et al., 2013), while sub-grid variabilities might be offset by averaging over large
455 watersheds. Therefore, we plotted the uncertainty against watershed areas in **Figure 7**, which verifies
456 that for most variables, the uncertainty declines as the watershed area increases. In addition to the general
457 understanding, **Figure 7** also reveals some interesting patterns. The most obvious decline comes from
458 ET (green) and longwave radiation (red), both of which are highly dependent on the land surface
459 conditions and are significantly affected by land surface spatial heterogeneity, thus benefiting the most
460 from spatial averaging for large river basins. Shortwave radiation and precipitation uncertainty show a
461 similar decline pattern (blue and purple), which is possibly related to their strong ties to cloud covers.
462 Temperature has a low uncertainty, and its relationship to watershed area is also not obvious. Wind speed
463 uncertainty only declines slightly as area increases, and we believe this is because wind speed uncertainty
464 can be traced back more to the atmospheric circulation patterns instead of the land surface conditions,
465 thus showing non-prominent relationship with watershed area. Overall, GSHA provides uncertainty
466 estimates that capture these prominent patterns, which can be helpful to hydrologic modellers and users.

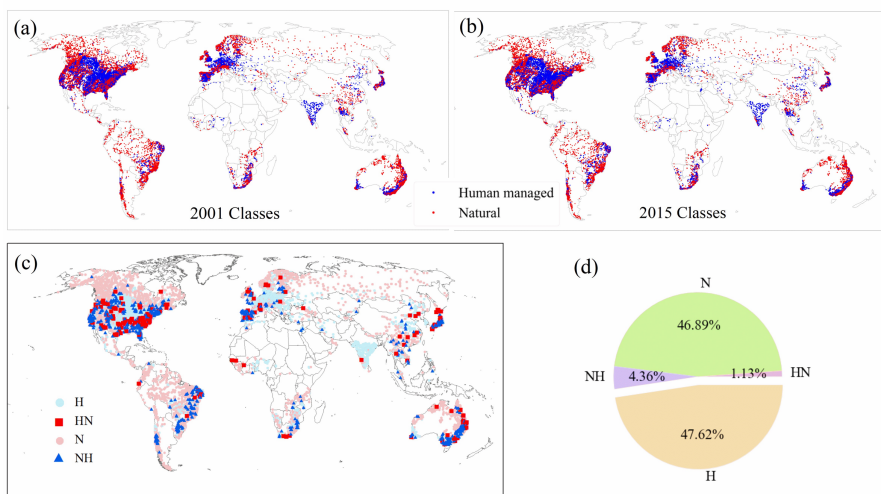
467 4.3 Natural and human managed watersheds and changing patterns

468 We also demonstrate the other key features of GSHA by categorizing global watersheds into
469 natural and human managed, and more prominently their temporal shifts in **Figure 8**. Overall,
470 majority of human managed watersheds locates in the US, Europe, and other regions with intensive
471 industrial or agricultural activities such as East and South Asia (**Figures 8a and 8b**). During 2001-
472 2015, 46.89% watersheds remained natural, while another 47.62% under human management in
473 2001 remained in the category throughout the study period (**Figure 8d**). Generally, northern
474 hemisphere has a larger proportion of human managed watersheds, while watersheds in the less
475 populated and urbanized southern hemisphere largely remain natural.

476 Noticeably, 4.36% of GSHA watersheds switched from natural to human managed (1011
477 watersheds), and the remaining 1.13% changed backed to natural states from human managed. For
478 instance, watersheds in the middle and lower Yangtze River area in China show a shift from human
479 managed to natural state, where environmental projects for ecological restoration were in place (Qu
480 et al., 2018; Zhang et al., 2015). Although the time span of GSHA LULC dynamics restricted the
481 change detection for developed regions as their urbanizations and infrastructure developments have
482 long been completed, and for fast emerging economies after 2015, the time series are also missing;
483 nevertheless, the changing human activities captured by GSHA may be helpful to understand the
484 streamflow changes including flood characteristics (Long Yang et al., 2021; Zhang et al., 2022).



Manuscript in submission to *ESSD*



485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

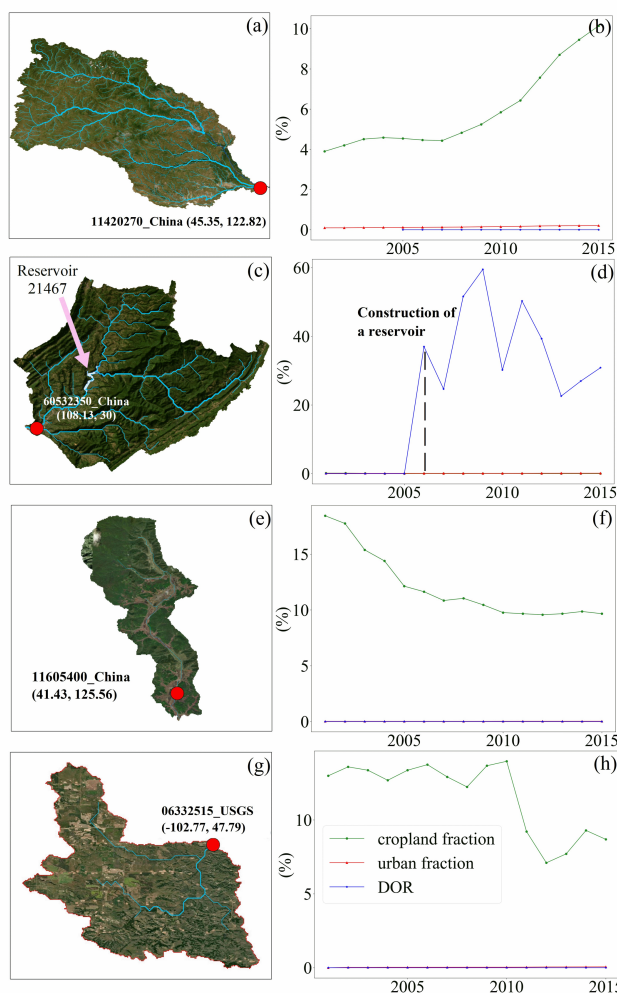
502

Figure 8 Classification of natural and human managed watersheds in 2001 (a) and 2015 (b). Changes in watershed categories are illustrated by (c) and (d). H and N in (c) and (d) represent watersheds that maintained human managed or natural from 2001-2015; NH and HN represent those changing from natural to human managed and from human managed to natural, respectively.

We further use several examples to illustrate the changing status of GSHA watersheds (**Figure 9**). **Figures 9a and 9b** show a watershed located in Northeastern China, where the rapid increase in cropland shifted the watershed from natural states to human managed in recent years. **Figures 9c and 9d** correspond to a mountainous area in Sichuan Province, China, which became human managed due to the construction of a reservoir in 2006. For another case in Northeastern China (**Figures 9e and 9f**) and a USGS case (**Figures 9g and 9h**), the watersheds shifted from human managed to natural, which is mainly manifested by the reduction in cropland fraction due to the environmental policy. For instance, afforestation during 2000-2010 in Changbai Mountains where the watershed in **Figures 9e and 9f** is located, significantly increased the forest cover and might bring a decline in human disturbance in the form of land use (Zhang & Liang, 2014). These results highlight the shifting watershed status that would require further attention from LSH users, which is encapsulated in GSHA v1.0 and will be continuously improved in the future.



Manuscript in submission to *ESSD*



503
 504 **Figure 9** Cases for shifting status of the watershed classification. (a) and (b) correspond to
 505 11420270_China, and (c) and (d) correspond to 60532350_China, both of which changed from natural
 506 to human managed category. (e) and (f) represent 11605400_China, and (g) and (h) correspond to
 507 06332515_USGS watershed changing from human managed to natural watershed.

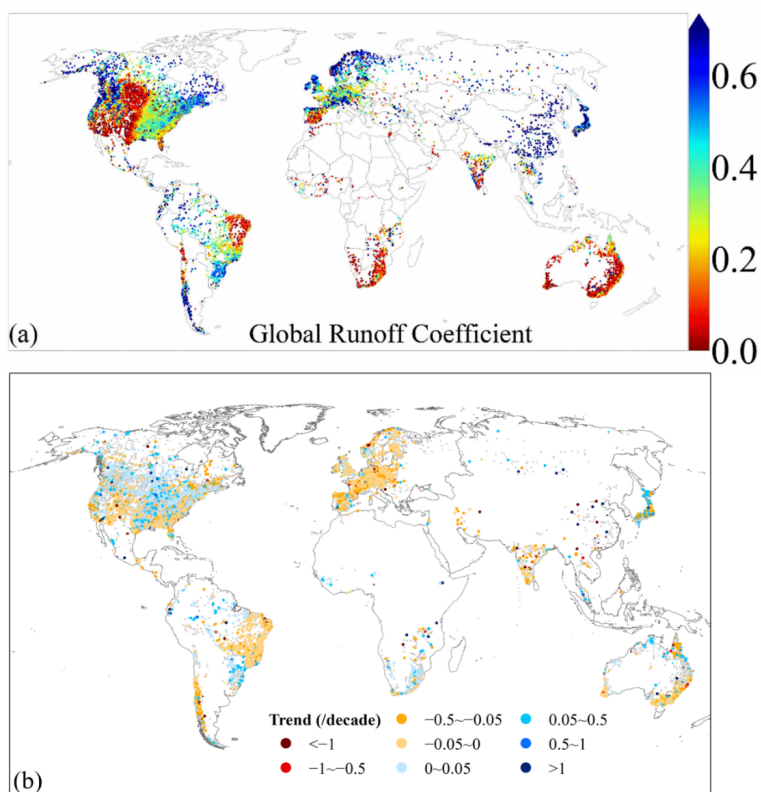
508 4.4 Changing runoff coefficient patterns derived from GSHA

509 Finally, we also analysed the spatial pattern and trend of global runoff coefficient (RC), as a
 510 brief demonstration of what GSHA can offer out of its many potential usages. RC is defined as R/P ,
 511 where R denotes runoff (mm) and P denotes precipitation (mm) in the corresponding period. **Figure**
 512 **10a** shows that, regions where a large proportion of rainfall goes into rivers instead of being
 513 evaporated or consumed (i.e., high RC) are in east Asia and North America, most parts of Europe,
 514 west coast of North America and the Amazon, in general agreement with the aridity patterns across
 515 the globe. For arid/semiarid areas and places with intense water use (e.g., western US, eastern Brazil,



Manuscript in submission to *ESSD*

516 Australia, Africa), RC is low, meaning most of the precipitation does not reach the gauged river.
517 Interestingly, we show the RC trend for the past decades. We found that RC generally remained
518 stable in most parts of the world (i.e., grey dots, **Figure 10b**), where >80% of the gauges do not
519 show statistically significant trend ($p > 0.05$). In total, 4252 watersheds demonstrate a 95%
520 significant trend in RC (5690 watersheds at a 90% significance level), and among them decreasing
521 RC is more widespread compared to increasing RC. The most pronounced RC decreasing trends are
522 observed in Europe, India, eastern Brazil, Chile, eastern Australia, and the Euphrates and Tigris,
523 which largely correspond to regions with known increasing agricultural, industrial, and residential
524 water use that may have reduced the river water. We also found that the global RC trend patterns
525 are different from a recent study showing mostly increasing RC in high-latitude watersheds, central
526 North America, eastern Australia, and Europe, which can largely be explained by the ET changes
527 (Xiong et al., 2022). Here GSHA reveals that RC trend is decreasing more widely than the increasing
528 trend, which may be concerning for the water availability in a changing climate, and will warrant
529 more future research along this line.



530
531 **Figure 10** Patterns of runoff coefficient (a) and its trend (b). Only watersheds with statistically
532 significant trend ($p < 0.05$) are shown with colours in (b); the small and large sized points represent 95%
533 ($p < 0.05$) and 90% significance level ($p < 0.1$), respectively. Note that the temporal coverage is different
534 for different gauges; readers can refer to the GSHA temporal coverage for interpreting the patterns. The



Manuscript in submission to *ESSD*

535 figure illustrates 18987 GSHA watersheds. Watersheds with less than 10 years of indices calculated
536 from over 250 valid observations per year, as well as with runoff coefficient trend over 20 per decade,
537 are not demonstrated in subfigure b.

538 5 Conclusions

539 Large sample hydrology (LSH) datasets play a critical role in data-driven analyses and model
540 parameter estimation for hydrological studies. From MOPEX (Duan et al., 2006) to Caravan
541 (Kratzert et al., 2023), significant efforts have been made to improve the comprehensiveness of LSH,
542 yet issues related to uncertainty estimates, and the human activity dynamics and the data spatial
543 coverage remain to be solved. This study focuses on complementing existing LSH with a new
544 synthesis dataset named the Global Streamflow characteristics, Hydrometeorology, and catchment
545 Attributes for large sample river-centric studies (GSHA v1.0).

546 To summarize, GSHA contributes the following aspects to the LSH development:

- 547 1. It includes streamflow indices, hydrometeorological data, and surface characteristics data for
548 21568 gauges compiled from 13 agencies worldwide, which represents one of the most
549 comprehensive LSH by far.
- 550 2. We incorporate multiple data sources to provide uncertainty estimates for each meteorological
551 variable (including precipitation, 2 m air temperature, radiation, wind, and ET). The spatial
552 patterns and the relationship between the uncertainty and the watershed characteristics that
553 GSHA revealed may be helpful to identify inconsistencies among data-driven studies or biases
554 for model parameter estimation studies using existing LSH.
- 555 3. Dynamic data are provided for previously static data descriptors for land cover changes
556 including urban, cropland and forest fractions, as well as reservoir storage change including
557 storage capacity and degree of regulation.

558 Although GSHA does not cover watersheds of $<25\text{km}^2$ or the dynamics of cryosphere variables
559 (e.g., glacier and permafrost) that become increasingly important in terrestrial hydrological changes,
560 and the time spans for the dynamic descriptors of LULC are unable to cover the critical periods for
561 the advanced and less-advanced economies due to the constraints with existing LULC data, GSHA
562 is utilized to unravel the following insights:

- 563 1. The uncertainty patterns vary between variables and geographical regions, indicating that the
564 interpretation of model and analysis results need to consider inconsistencies of raw data, apart
565 from looking into the methodologies and patterns themselves.
- 566 2. Although most watersheds have remained natural or human managed throughout the GSHA
567 time span, a considerable number of watersheds shifted between the two categories, which can
568 be ascribed to urbanization, cropland increase, reservoir construction and ecological restoration
569 such as returning farmland to natural states, and these can be clearly manifested using GSHA.
- 570 3. Analysis with runoff coefficient reveals that while ~80% of gauges do not observe a
571 statistically significant trend, a greater portion of gauges have experienced a declining RC trend
572 than an increase trend. This pattern revealed by GSHA can be used to further study water
573 availability issues in a changing climate.



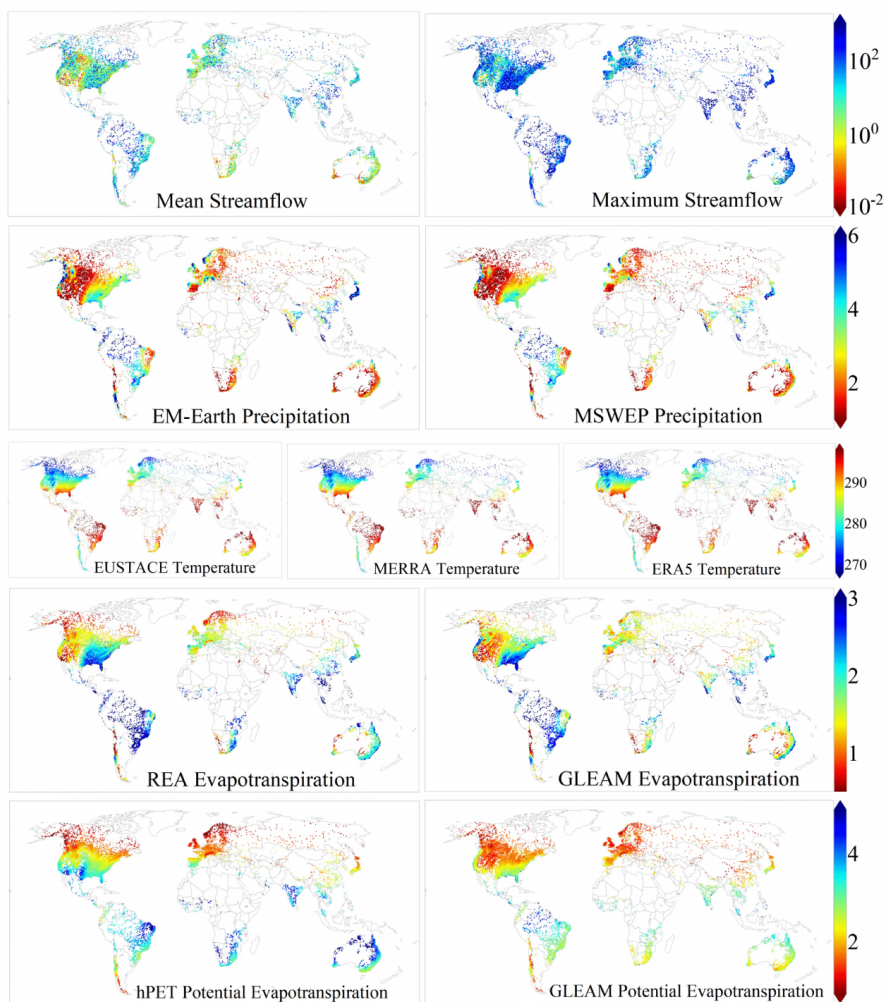
Manuscript in submission to *ESSD*

574 As our knowledge on the above processes continues to improve, we expect that future versions
575 of GSHA will be continuously updated. Finally, better hydrological data sharing is crucial to
576 advance global change hydrology studies.

577 Appendix

578 A. Spatial patterns of GSHA meteorological variables

579 **Figures A1 & A2** show the spatial distributions of GSHA meteorological variables and selected
580 streamflow indices. The spatial pattern derived from each individual data source is plotted separately.

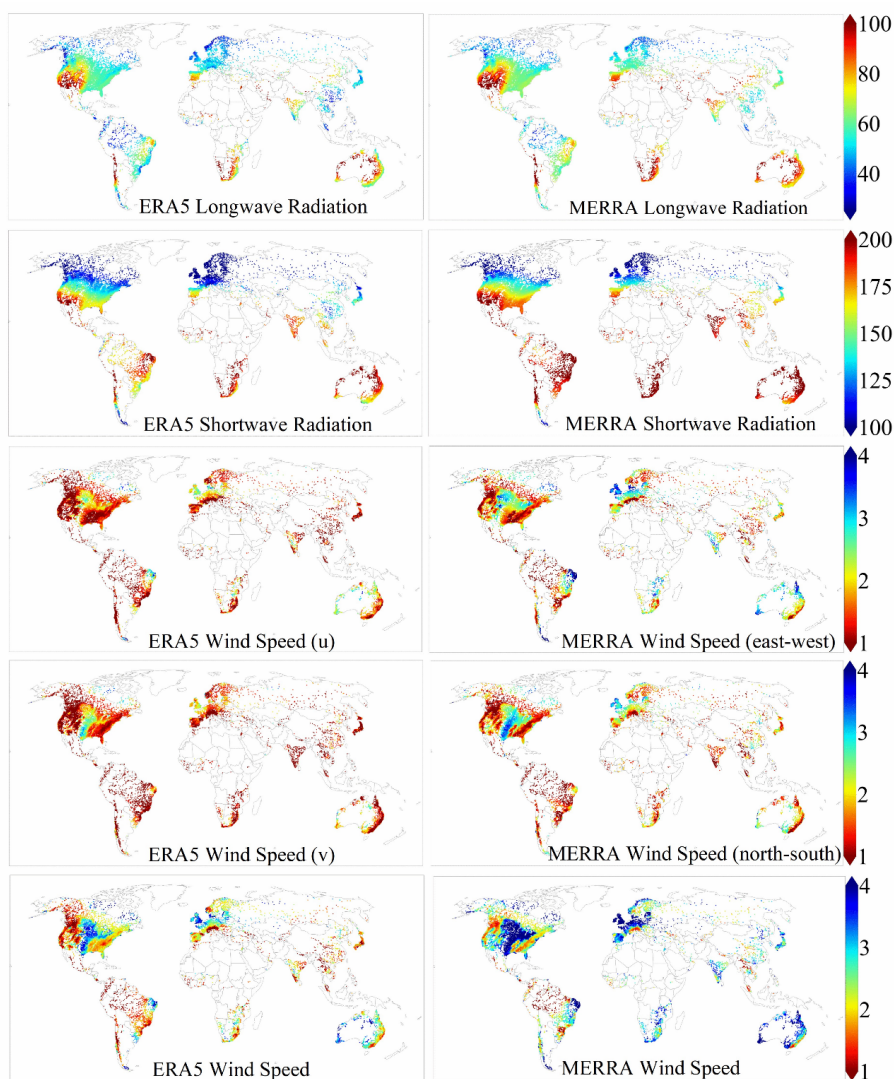


581
582 **Figure A1** Spatial distribution of streamflow indices (row 1, m^3/s), precipitation (row 2, mm/day), 2 m



Manuscript in submission to *ESSD*

583 air temperature (row 3, K), actual ET (row 4, mm/day), potential ET (row 5, mm/day).



584

585 **Figure A 2** Spatial distribution of longwave radiation (row 1, W/m²), shortwave radiation (row 2, W/m²),
586 wind u- (row 3, m/s) and v- components (row 4, m/s) and the wind speed (row 5, m/s).

587

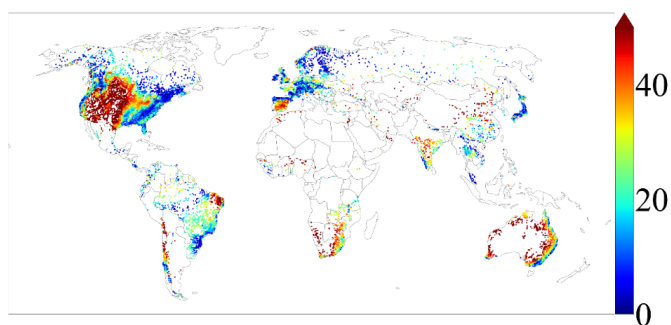
588 **B. Potential evapotranspiration uncertainty**

589 The spatial and numerical distributions of potential evapotranspiration (PET) uncertainties are
590 illustrated in **Figure B1** and **Figure B2**. PET uncertainty is high compared with other variables (see
591 5.2 section). The majority of high PET uncertainty watersheds are in dry areas, but since it is
592 calculated from meteorological variables, exceptions exist for places including eastern Pacific coast,



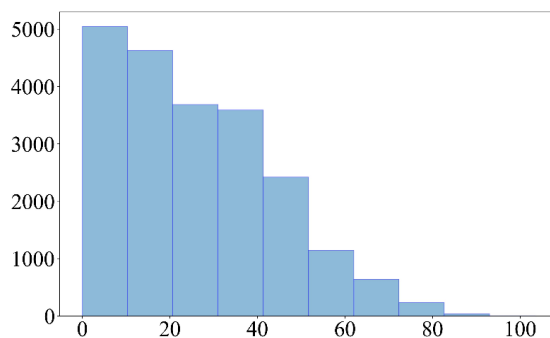
Manuscript in submission to *ESSD*

593 where the climate is dry but PET uncertainty is low, and India, which is located in a wet climate
594 zone but has high PET uncertainty. As demonstrated by **Figure B3**, PET uncertainty do not decrease
595 with the increase of watershed area, probably because PET is calculated from various variables, and
596 the calculation over large watersheds involves more uncertainties for individual grids.



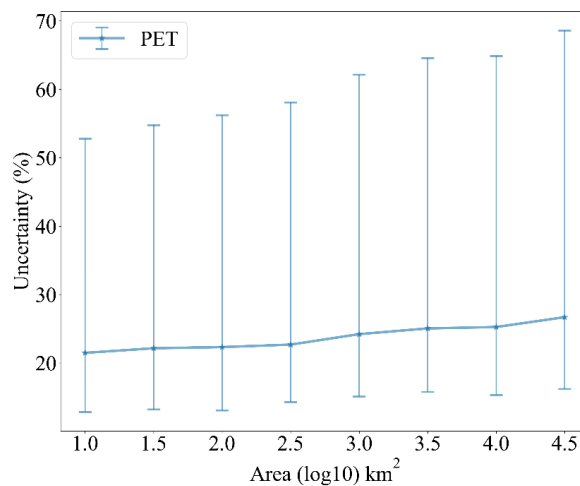
597
598
599

Figure B1 Spatial pattern of potential evapotranspiration (PET) uncertainty.



600
601
602

Figure B2 Numerical distribution of PET uncertainty.



603
604

Figure B3 Relationship of PET uncertainty to watershed area.



Manuscript in submission to *ESSD*

605 Author contribution

606 Conceptualization: PL. Investigation: ZY, PL, RR, GA, XL. Data curation: ZY, RR, XL, PL, ZZ,
607 SC. Funding acquisition: PL. Writing - initial: ZY, PL. Writing - Review and Editing: PL, ZY, GA,
608 RR, XL.

609 Data and Code Availability

610 GSHA v1.0 is openly available at <https://doi.org/10.5281/zenodo.8090704> (Yin et al. 2023). The
611 codes involved in the workflow to generating GSHA will be available upon reasonable requests to
612 the corresponding author.

613 Competing interests

614 The authors declare no conflict of interest.

615 Acknowledgements

616 This study is supported by the National Key Research and Development Program
617 (2022YFF0801303), and the Fundamental Research Funds for the Central Universities, Peking
618 University on “Numerical modelling and remote sensing of global river discharge” (#7100604136).

619 References

- 620 Addor, N., Do, H. X., Alvarez-Garreton, C., Coxon, G., Fowler, K., & Mendoza, P. A. (2020). Large-
621 sample hydrology: recent progress, guidelines for new datasets and grand challenges.
622 *Hydrological Sciences Journal-Journal Des Sciences Hydrologiques*, 65(5), 712-725.
623 <https://doi.org/10.1080/02626667.2019.1683182>
- 624 Addor, N., Nearing, G., Prieto, C., Newman, A. J., Le Vine, N., & Clark, M. P. (2018). A ranking of
625 hydrological signatures based on their predictability in space. *Water Resources Research*, 54(11),
626 8792-8812.
- 627 Addor, N., Newman, A. J., Mizukami, N., & Clark, M. P. (2017). The CAMELS data set: catchment
628 attributes and meteorology for large-sample studies. *Hydrology and Earth System Sciences*,
629 21(10), 5293-5313. <https://doi.org/10.5194/hess-21-5293-2017>
- 630 Aerts, J. P., Hut, R. W., van de Giesen, N. C., Drost, N., van Verseveld, W. J., Weerts, A. H., & Hazenberg,
631 P. (2022). Large-sample assessment of varying spatial resolution on the streamflow estimates of
632 the wflow_sbm hydrological model. *Hydrology and Earth System Sciences*, 26(16), 4407-4430.



Manuscript in submission to *ESSD*

- 633 AghaKouchak, A., Chiang, F., Huning, L. S., Love, C. A., Mallakpour, I., Mazdiyasi, O., Moftakhari,
634 H., Papalexioiu, S. M., Ragno, E., & Sadegh, M. (2020). Climate Extremes and Compound
635 Hazards in a Warming World. *Annual Review of Earth and Planetary Sciences, Vol 48, 2020*,
636 48, 519-548. <Go to ISI>://WOS:000613951000021
- 637 Alvarez-Garreton, C., Mendoza, P. A., Boisier, J. P., Addor, N., Galleguillos, M., Zambrano-Bigiarini,
638 M., Lara, A., Puelma, C., Cortes, G., Garreaud, R., McPhee, J., & Ayala, A. (2018). The
639 CAMELS-CL dataset: catchment attributes and meteorology for large sample studies - Chile
640 dataset. *Hydrology and Earth System Sciences*, 22(11), 5817-5846.
641 <https://doi.org/10.5194/hess-22-5817-2018>
- 642 ArcticNET 2022 ArcticNET V1.0 (available at: <https://russia-arcticnet.sr.unh.edu/>)
- 643 Arsenault, R., Brissette, F., Martel, J.-L., Troin, M., Lévesque, G., Davidson-Chaput, J., Gonzalez, M.
644 C., Ameli, A., & Poulin, A. (2020). A comprehensive, multisource database for
645 hydrometeorological modeling of 14,425 North American watersheds. *Scientific Data*, 7(1), 243.
- 646 Australian Bureau of Meteorology 2022 Australian Bureau of Meteorology (available at:
647 www.bom.gov.au/)
- 648 Beck, H. E., van Dijk, A. I., De Roo, A., Miralles, D. G., McVicar, T. R., Schellekens, J., & Bruijnzeel,
649 L. A. (2016). Global-scale regionalization of hydrologic model parameters. *Water Resources*
650 *Research*, 52(5), 3599-3622.
- 651 Beck, H. E., Van Dijk, A. I., Levizzani, V., Schellekens, J., Miralles, D. G., Martens, B., & De Roo, A.
652 (2017). MSWEP: 3-hourly 0.25 global gridded precipitation (1979–2015) by merging gauge,
653 satellite, and reanalysis data. *Hydrology and Earth System Sciences*, 21(1), 589-615.
- 654 Beck, H. E., Wood, E. F., Pan, M., Fisher, C. K., Miralles, D. G., Van Dijk, A. I., McVicar, T. R., & Adler,
655 R. F. (2019). MSWEP V2 global 3-hourly 0.1 precipitation: methodology and quantitative
656 assessment. *Bulletin of the American Meteorological Society*, 100(3), 473-500.
- 657 Belvederesi, C., Zaghoul, M. S., Achari, G., Gupta, A., & Hassan, Q. K. (2022). Modelling river flow in
658 cold and ungauged regions: A review of the purposes, methods, and challenges. *Environmental*
659 *Reviews*, 30(1), 159-173.
- 660 Benke, K. K., Lowell, K. E., & Hamilton, A. J. (2008). Parameter uncertainty, sensitivity analysis and
661 prediction error in a water-balance hydrological model. *Mathematical and Computer Modelling*,
662 47(11-12), 1134-1149.
- 663 Beven, K. J., & Alcock, R. E. (2012). Modelling everything everywhere: a new approach to decision-
664 making for water management under uncertainty. *Freshwater Biology*, 57, 124-132.
- 665 Bourdin, D. R., Fleming, S. W., & Stull, R. B. (2012). Streamflow modelling: a primer on applications,
666 approaches and challenges. *Atmosphere-Ocean*, 50(4), 507-536.
- 667 Brazil National Water Agency 2022 National water and sanitation agency (ANA) Agência Nac Águas E
668 Saneam. Básico ANA (available at: www.gov.br/ana/en/national_water_agency)
- 669 Brugnara, Y., Good, E., Squintu, A. A., van der Schrier, G., & Brönnimann, S. (2019). The EUSTACE
670 global land station daily air temperature dataset. *Geoscience Data Journal*, 6(2), 189-204.
- 671 Brun, P., Zimmermann, N. E., Hari, C., Pellissier, L., & Karger, D. N. (2022). Global climate-related
672 predictors at kilometer resolution for the past and future. *Earth System Science Data*, 14(12),
673 5573-5603.
- 674 Brunner, M. I., Slater, L., Tallaksen, L. M., & Clark, M. (2021). Challenges in modeling and predicting
675 floods and droughts: A review. *Wiley Interdisciplinary Reviews: Water*, 8(3), e1520.
- 676 Burges, S. J. (1998). Streamflow prediction: capabilities, opportunities, and challenges. *Hydrologic*



Manuscript in submission to *ESSD*

- 677 *Sciences: Taking Stock and Looking Ahead*, 5, 101-134.
- 678 Canada National Water Data Archive 2022 National water data archive HYDAT (available at:
679 [www.canada.ca/en/environment-climate-change/services/water-](http://www.canada.ca/en/environment-climate-change/services/water-overview/quantity/monitoring/survey/data-products-services/national-archive-hydat.html)
680 [overview/quantity/monitoring/survey/data-products-services/national-archive-hydat.html](http://www.canada.ca/en/environment-climate-change/services/water-overview/quantity/monitoring/survey/data-products-services/national-archive-hydat.html))
- 681 Chagas, V. B., Chaffe, P. L., Addor, N., Fan, F. M., Fleischmann, A. S., Paiva, R. C., & Siqueira, V. A.
682 (2020). CAMELS-BR: hydrometeorological time series and landscape attributes for 897
683 catchments in Brazil. *Earth System Science Data*, 12(3), 2075-2096.
- 684 Chen, X., Jiang, L., Luo, Y., & Liu, J. (2023). A global streamflow indices time series dataset for large-
685 sample hydrological analyses on streamflow regime (until 2021). *Earth System Science Data*
686 *Discussions*, 2023, 1-18.
- 687 Chile Center for Climate and Resilience Research 2022 Center for climate and resilience research CR2 |
688 Chilean research center on climate, climate change and resilience (available at: www.cr2.cl/eng/)
- 689 Cho, K., & Kim, Y. (2022). Improving streamflow prediction in the WRF-Hydro model with LSTM
690 networks. *Journal of Hydrology*, 605, 127297.
- 691 Clark, M. P., Vogel, R. M., Lamontagne, J. R., Mizukami, N., Knoben, W. J., Tang, G., Gharari, S., Freer,
692 J. E., Whitfield, P. H., & Shook, K. R. (2021). The abuse of popular performance metrics in
693 hydrologic modeling. *Water Resources Research*, 57(9), e2020WR029001.
- 694 Claverie, M., Matthews, J. L., Vermote, E. F., & Justice, C. O. (2016). A 30+ year AVHRR LAI and
695 FAPAR climate data record: Algorithm description and validation. *Remote Sensing*, 8(3), 263.
- 696 Coxon, G., Addor, N., Bloomfield, J. P., Freer, J., Fry, M., Hannaford, J., Howden, N. J., Lane, R., Lewis,
697 M., & Robinson, E. L. (2020). CAMELS-GB: hydrometeorological time series and landscape
698 attributes for 671 catchments in Great Britain. *Earth System Science Data*, 12(4), 2459-2483.
- 699 Delaigue, O., Brigode, P., Andréassian, V., Perrin, C., Etchevers, P., Soubeyrou, J.-M., Janet, B., & Nans,
700 A. (2022). CAMELS-FR: A large sample hydroclimatic dataset for France to explore
701 hydrological diversity and support model benchmarking. *IAHS-2022 Scientific Assembly*, 575.
- 702 Do, H. X., Gudmundsson, L., Leonard, M., & Westra, S. (2018). The Global Streamflow Indices and
703 Metadata Archive (GSIM) - Part 1: The production of a daily streamflow archive and metadata.
704 *Earth System Science Data*, 10(2). <https://doi.org/10.5194/essd-10-765-2018>
- 705 Duan, Q., Schaake, J., Andréassian, V., Franks, S., Goteti, G., Gupta, H., Gusev, Y., Habets, F., Hall, A.,
706 & Hay, L. (2006). Model Parameter Estimation Experiment (MOPEX): An overview of science
707 strategy and major results from the second and third workshops. *Journal of Hydrology*, 320(1-
708 2), 3-17.
- 709 Fang, Y., Huang, Y., Qu, B., Zhang, X., Zhang, T., & Xia, D. (2022). Estimating the Routing Parameter
710 of the Xin'anjiang Hydrological Model Based on Remote Sensing Data and Machine Learning.
711 *Remote Sensing*, 14(18), 4609.
- 712 Fowler, K. J. A., Acharya, S. C., Addor, N., Chou, C. C., & Peel, M. C. (2021). CAMELS-AUS:
713 hydrometeorological time series and landscape attributes for 222 catchments in Australia. *Earth*
714 *System Science Data*, 13(8), 3847-3867. <https://doi.org/10.5194/essd-13-3847-2021>
- 715 Friedl, M. A., Sulla-Menashe, D., Tan, B., Schneider, A., Ramankutty, N., Sibley, A., & Huang, X. (2010).
716 MODIS Collection 5 global land cover: Algorithm refinements and characterization of new
717 datasets. *Remote Sensing of Environment*, 114(1), 168-182.
- 718 Friedl, M., D. Sulla-Menashe. MCD12Q1 MODIS/Terra+Aqua Land Cover Type Yearly L3 Global 500m
719 SIN Grid V006. 2019, distributed by NASA EOSDIS Land Processes DAAC,
720 <https://doi.org/10.5067/MODIS/MCD12Q1.006>.



Manuscript in submission to *ESSD*

- 721 Gelaro, R., McCarty, W., Suárez, M. J., Todling, R., Molod, A., Takacs, L., Randles, C. A., Darmenov,
722 A., Bosilovich, M. G., & Reichle, R. (2017). The modern-era retrospective analysis for research
723 and applications, version 2 (MERRA-2). *Journal of Climate*, 30(14), 5419-5454.
- 724 Global Modeling and Assimilation Office (GMAO) (2015), inst3_3d_asm_Cp: MERRA-2 3D IAU State,
725 Meteorology Instantaneous 3-hourly (p-coord, 0.625x0.5L42), version 5.12.4, Greenbelt, MD,
726 USA: Goddard Space Flight Center Distributed Active Archive Center (GSFC DAAC), doi:
727 10.5067/VJAFPLIICSIV.
- 728 Gudmundsson, L., Do, H. X., Leonard, M., & Westra, S. (2018). The Global Streamflow Indices and
729 Metadata Archive (GSIM) - Part 2: Quality control, time-series indices and homogeneity
730 assessment. *Earth System Science Data*, 10(2). <https://doi.org/10.5194/essd-10-787-2018>
- 731 Gupta, H. V., Perrin, C., Blöschl, G., Montanari, A., Kumar, R., Clark, M., & Andreassian, V. (2014).
732 Large-sample hydrology: a need to balance depth with breadth. *Hydrology and Earth System
733 Sciences*, 18(2), 463-477. <https://doi.org/10.5194/hess-18-463-2014>
- 734 Hao, Z., Jin, J., Xia, R., Tian, S., Yang, W., Liu, Q., Zhu, M., Ma, T., Jing, C., & Zhang, Y. (2021). CCAM:
735 China catchment attributes and meteorology dataset. *Earth System Science Data*, 13(12), 5591-
736 5616.
- 737 Henck, A. C., Montgomery, D. R., Huntington, K. W., & Liang, C. (2010). Monsoon control of effective
738 discharge, Yunnan and Tibet. *Geology*, 38(11), 975-978.
- 739 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., Nicolas, J., Peubey,
740 C., Radu, R., & Schepers, D. (2020). The ERA5 global reanalysis. *Quarterly Journal of the
741 Royal Meteorological Society*, 146(730), 1999-2049.
- 742 Hrachowitz, M., Savenije, H., Blöschl, G., McDonnell, J., Sivapalan, M., Pomeroy, J., Arheimer, B.,
743 Blume, T., Clark, M., & Ehret, U. (2013). A decade of Predictions in Ungauged Basins (PUB)—
744 a review. *Hydrological sciences journal*, 58(6), 1198-1255.
- 745 Hu, D., Cao, S., Chen, S., Deng, L., & Feng, N. (2017). Monitoring spatial patterns and changes of
746 surface net radiation in urban and suburban areas using satellite remote-sensing data.
747 *International Journal of Remote Sensing*, 38(4), 1043-1061.
- 748 Huang, Yinghui (2020): High spatiotemporal resolution mapping of global urban change from 1985 to
749 2015. figshare. Dataset. <https://doi.org/10.6084/m9.figshare.11513178.v1>
- 750 Immerzeel, W., and, & Droogers, P. (2008). Calibration of a distributed hydrological model based on
751 satellite evapotranspiration. *Journal of Hydrology*, 349(3-4), 411-424.
- 752 India Water Resources Information System 2022 India Water Resources Information System (available
753 at: <https://indiawris.gov.in/wris/#/>)
- 754 Japanese Water Information System 2022 Ministry of Land, Infrastructure, Transport and Tourism
755 (available at: www.mlit.go.jp/en/)
- 756 Kauffeldt, A., Halldin, S., Rodhe, A., Xu, C.-Y., & Westerberg, I. K. (2013). Disinformative data in large-
757 scale hydrological modelling. *Hydrology and Earth System Sciences*, 17(7), 2845-2857.
- 758 Klingler, C., Schulz, K., & Herrnegger, M. (2021). LamaH-CE: LARge-SaMple DATA for Hydrology and
759 Environmental Sciences for Central Europe. *Earth System Science Data*, 13(9), 4529-4565.
760 <https://doi.org/10.5194/essd-13-4529-2021>
- 761 Kovács, G. (1984). Proposal to construct a coordinating matrix for comparative hydrology. *Hydrological
762 sciences journal*, 29(4), 435-443.
- 763 Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., & Nearing, G. (2019a). Benchmarking
764 a catchment-aware long short-term memory network (LSTM) for large-scale hydrological



Manuscript in submission to *ESSD*

- 765 modeling. *Hydrol. Earth Syst. Sci. Discuss*, 2019, 1-32.
- 766 Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., & Nearing, G. (2019b). Towards
767 learning universal, regional, and local hydrological behaviors via machine learning applied to
768 large-sample datasets. *Hydrology and Earth System Sciences*, 23(12), 5089-5110.
- 769 Kratzert, F., Nearing, G., Addor, N., Erickson, T., Gauch, M., Gilon, O., Gudmundsson, L., Hassidim, A.,
770 Klotz, D., & Nevo, S. (2023). Caravan-A global community dataset for large-sample hydrology.
771 *Scientific Data*, 10(1), 61.
- 772 Lehner, B., Messenger, M. L., Korver, M. C., & Linke, S. (2022). Global hydro-environmental lake
773 characteristics at high spatial resolution. *Scientific Data*, 9(1), 351.
- 774 Li, B., Rodell, M., Kumar, S., Beaudoin, H. K., Getirana, A., Zaitchik, B. F., de Goncalves, L. G.,
775 Cossetin, C., Bhanja, S., & Mukherjee, A. (2019). Global GRACE data assimilation for
776 groundwater and drought monitoring: Advances and challenges. *Water Resources Research*,
777 55(9), 7564-7586.
- 778 Lin, P., Rajib, M. A., Yang, Z. L., Somos-Valenzuela, M., Merwade, V., Maidment, D. R., Wang, Y., &
779 Chen, L. (2018). Spatiotemporal evaluation of simulated evapotranspiration and streamflow
780 over Texas using the WRF-Hydro-RAPID modeling framework. *JAWRA Journal of the*
781 *American Water Resources Association*, 54(1), 40-54.
- 782 Lin, P. R., Pan, M., Beck, H. E., Yang, Y., Yamazaki, D., Frasson, R., David, C. H., Durand, M., Pavelsky,
783 T. M., Allen, G. H., Gleason, C. J., & Wood, E. F. (2019). Global Reconstruction of Naturalized
784 River Flows at 2.94 Million Reaches. *Water Resources Research*, 55(8), 6499-6516.
785 <https://doi.org/10.1029/2019wr025287>
- 786 Lin, P. R., Pan, M., Wood, E. F., Yamazaki, D., & Allen, G. H. (2021). A new vector-based global river
787 network dataset accounting for variable drainage density. *Scientific Data*, 8(1).
788 <https://doi.org/ARTN2810.1038/s41597-021-00819-9>
- 789 Linke, S., Lehner, B., Ouellet Dallaire, C., Ariwi, J., Grill, G., Anand, M., Beames, P., Burchard-Levine,
790 V., Maxwell, S., & Moidu, H. (2019). Global hydro-environmental sub-basin and river reach
791 characteristics at high spatial resolution. *Scientific Data*, 6(1), 283.
- 792 Liu, X., Huang, Y., Xu, X., Li, X., Li, X., Ciais, P., Lin, P., Gong, K., Ziegler, A. D., & Chen, A. (2020).
793 High-spatiotemporal-resolution mapping of global urban change from 1985 to 2015. *Nature*
794 *Sustainability*, 3(7), 564-570.
- 795 Lu, J., Wang, G., Chen, T., Li, S., Hagan, D. F. T., Kattel, G., Peng, J., Jiang, T., & Su, B. (2021). A
796 harmonized global land evaporation dataset from model-based products covering 1980–2017.
797 *Earth System Science Data*, 13(12), 5879-5898.
- 798 Martens, B., Miralles, D. G., Lievens, H., Van Der Schalie, R., De Jeu, R. A., Fernández-Prieto, D., Beck,
799 H. E., Dorigo, W. A., & Verhoest, N. E. (2017). GLEAM v3: Satellite-based land evaporation
800 and root-zone soil moisture. *Geoscientific Model Development*, 10(5), 1903-1925.
- 801 Merchant, C. J., Paul, F., Popp, T., Ablain, M., Bontemps, S., Defourny, P., Hollmann, R., Lavergne, T.,
802 Laeng, A., & De Leeuw, G. (2017). Uncertainty information in climate data records from Earth
803 observation. *Earth System Science Data*, 9(2), 511-527.
- 804 Miralles, D. G., Holmes, T., De Jeu, R., Gash, J., Meesters, A., & Dolman, A. (2011). Global land-surface
805 evaporation estimated from satellite-based observations. *Hydrology and Earth System Sciences*,
806 15(2), 453-469.
- 807 Muñoz-Sabater, J., Dutra, E., Agustí-Panareda, A., Albergel, C., Arduini, G., Balsamo, G., Boussetta, S.,
808 Choulga, M., Harrigan, S., & Hersbach, H. (2021). ERA5-Land: A state-of-the-art global



Manuscript in submission to *ESSD*

- 809 reanalysis dataset for land applications. *Earth System Science Data*, 13(9), 4349-4383.
- 810 Muñoz Sabater, J., (2019): ERA5-Land hourly data from 1981 to present. Copernicus Climate Change
811 Service (C3S) Climate Data Store (CDS), 10.24381/cds.e2161bac
- 812 Nandi, S., & Reddy, M. J. (2022). An integrated approach to streamflow estimation and flood inundation
813 mapping using VIC, RAPID and LISFLOOD-FP. *Journal of Hydrology*, 610, 127842.
- 814 Newman, A. J., Clark, M. P., Sampson, K., Wood, A., Hay, L. E., Bock, A., Viger, R. J., Blodgett, D.,
815 Brekke, L., Arnold, J. R., Hopson, T., & Duan, Q. (2015). Development of a large-sample
816 watershed-scale hydrometeorological data set for the contiguous USA: data set characteristics
817 and assessment of regional variability in hydrologic model performance. *Hydrology and Earth
818 System Sciences*, 19(1), 209-223. <https://doi.org/10.5194/hess-19-209-2015>
- 819 Niraula, R., Meixner, T., & Norman, L. M. (2015). Determining the importance of model calibration for
820 forecasting absolute/relative changes in streamflow from LULC and climate changes. *Journal
821 of Hydrology*, 522, 439-451.
- 822 Qu, S., Wang, L., Lin, A., Zhu, H., & Yuan, M. (2018). What drives the vegetation restoration in Yangtze
823 River basin, China: climate change or anthropogenic factors? *Ecological Indicators*, 90, 438-
824 450.
- 825 Razavi, T., & Coulibaly, P. (2013). Streamflow prediction in ungauged basins: review of regionalization
826 methods. *Journal of hydrologic engineering*, 18(8), 958-975.
- 827 Ren, K., Fang, W., Qu, J., Zhang, X., & Shi, X. (2020). Comparison of eight filter-based feature selection
828 methods for monthly streamflow forecasting—three case studies on CAMELS data sets. *Journal
829 of Hydrology*, 586, 124897.
- 830 Riggs, R. M., Allen, G. H., Wang, J., Pavelsky, T. M., Gleason, C. J., David, C. H., & Durand, M. (2023).
831 Extending global river gauge records using satellite observations. *Environmental Research
832 Letters*.
- 833 Schaake, J., Cong, S., & Duan, Q. (2006). *US MOPEX data set*.
- 834 Schmidt, A. H., Montgomery, D. R., Huntington, K. W., & Liang, C. (2011). The question of communist
835 land degradation: new evidence from local erosion and basin-wide sediment yield in Southwest
836 China and Southeast Tibet. *Annals of the Association of American Geographers*, 101(3), 477-
837 496.
- 838 Schreiner-McGraw, A. P., & Ajami, H. (2020). Impact of uncertainty in precipitation forcing data sets on
839 the hydrologic budget of an integrated hydrologic model in mountainous terrain. *Water
840 Resources Research*, 56(12), e2020WR027639.
- 841 Singer, M. B., Asfaw, D. T., Rosolem, R., Cuthbert, M. O., Miralles, D. G., MacLeod, D., Quichimbo, E.
842 A., & Michaelides, K. (2021). Hourly potential evapotranspiration at 0.1 resolution for the
843 global land surface from 1981-present. *Scientific Data*, 8(1), 224.
- 844 Spain Anuario de Aforos 2022 Anuario de Aforos—Anuario de Aforos Digital—datos.gob.es (available
845 at: <http://datos.gob.es/es/catalogo/e00125801-anuario-de-aforos/resource/4836b826-e7fd-4a41-950c-89b4eaea0279>)
846
- 847 Sörensson, A. A., & Ruscica, R. C. (2018). Intercomparison and uncertainty assessment of nine
848 evapotranspiration estimates over South America. *Water Resources Research*, 54(4), 2891-2908.
- 849 Tang, G., Clark, M. P., & Papalexiou, S. M. (2022). EM-Earth: The ensemble meteorological dataset for
850 planet Earth. *Bulletin of the American Meteorological Society*, 103(4), E996-E1018.
- 851 Tang, G., Clark, M., Papalexiou, S. (2022) EM-Earth: The Ensemble Meteorological Dataset for Planet
852 Earth. Federated Research Data Repository. <https://doi.org/10.20383/102.0547>



Manuscript in submission to *ESSD*

- 853 Tang, G., Clark, M. P., Knoben, W. J. M., Liu, H., Gharari, S., Arnal, L., et al. (2023). The impact of
854 meteorological forcing uncertainty on hydrological modeling: A global analysis of cryosphere
855 basins. *Water Resources Research*, 59, e2022WR033767. <https://doi.org/10.1029/2022WR033767>
- 856 Thackeray, C. W., Hall, A., Norris, J., & Chen, D. (2022). Constraining the increased frequency of global
857 precipitation extremes under warming. *Nature Climate Change*, 12(5), 441-448.
858 <https://doi.org/10.1038/s41558-022-01329-1>
- 859 Thailand Royal Irrigation Department 2022 RID River Discharge Data (available at: [http://hydro.iis.u-](http://hydro.iis.u-tokyo.ac.jp/GAME-T/GAIN-T/routine/rid-river/disc_d.html)
860 [tokyo.ac.jp/GAME-T/GAIN-T/routine/rid-river/disc_d.html](http://hydro.iis.u-tokyo.ac.jp/GAME-T/GAIN-T/routine/rid-river/disc_d.html))
- 861 The Global Runoff Data Centre 2022 The global runoff data centre GRDC Data Portal (available at:
862 <https://portal.grdc.bafg.de/applications/public.html?publicuser=PublicUser>)
- 863 Ukhurebor, K. E., Azi, S. O., Aigbe, U. O., Onyancha, R. B., & Emegha, J. O. (2020). Analyzing the
864 uncertainties between reanalysis meteorological data and ground measured meteorological data.
865 *Measurement*, 165, 108110.
- 866 U.S. Geological Survey 2019 Gages Through the Ages (available at:
867 <https://labs.waterdata.usgs.gov/visualizations/gages-through-the-ages>)
- 868 Vermote, Eric; NOAA CDR Program. (2019): NOAA Climate Data Record (CDR) of AVHRR Leaf Area
869 Index (LAI) and Fraction of Absorbed Photosynthetically Active Radiation (FAPAR), Version
870 5. [LAI]. NOAA National Centers for Environmental Information.
871 <https://doi.org/10.7289/V5TT4P69>.
- 872 Wang, J., Walter, B. A., Yao, F., Song, C., Ding, M., Maroof, A. S., Zhu, J., Fan, C., McAlister, J. M., &
873 Sikder, S. (2022). GeoDAR: georeferenced global dams and reservoirs dataset for bridging
874 attributes and geolocations. *Earth System Science Data*, 14(4), 1869-1899.
- 875 Wilby, R., & Dessai, S. (2010). Robust adaptation to climate change.
- 876 Xiong, J., Yin, J., Guo, S., He, S., & Chen, J. (2022). Annual runoff coefficient variation in a changing
877 environment: A global perspective. *Environmental Research Letters*, 17(6), 064006.
- 878 Yamazaki, D., Ikeshima, D., Sosa, J., Bates, P. D., Allen, G. H., & Pavelsky, T. M. (2019). MERIT Hydro:
879 a high-resolution global hydrography map based on latest topography dataset. *Water Resources*
880 *Research*, 55(6), 5053-5073.
- 881 Yamazaki, D., Ikeshima, D., Tawatari, R., Yamaguchi, T., O'Loughlin, F., Neal, J. C., Sampson, C. C.,
882 Kanae, S., & Bates, P. D. (2017). A high-accuracy map of global terrain elevations. *Geophysical*
883 *Research Letters*, 44(11), 5844-5853.
- 884 Yang, L., Yang, Y., Villarini, G., Li, X., Hu, H., Wang, L., Blöschl, G., & Tian, F. (2021). Climate more
885 important for Chinese flood changes than reservoirs and land use. *Geophysical Research Letters*,
886 48(11), e2021GL093061.
- 887 Yang, L., Yang, Y. X., Villarini, G., Li, X., Hu, H. C., Wang, L. C., Blöschl, G., & Tian, F. Q. (2021).
888 Climate More Important for Chinese Flood Changes Than Reservoirs and Land Use.
889 *Geophysical Research Letters*, 48(11). <https://doi.org/ARTN>
890 e2021GL093061.1029/2021GL093061
- 891 Yin, Z., Lin, P., Riggs, R., Allen, G. H., Lei, X., Zheng, Z., & Cai, S. (2023). A Synthesis of Global
892 Streamflow characteristics, Hydrometeorology, and catchment Attributes (GSHA) for Large
893 Sample River-Centric Studies (1.0) [Data set]. Zenodo.
894 <https://doi.org/10.5281/zenodo.8090704>
- 895 Zaitchik, B. F., Rodell, M., & Reichle, R. H. (2008). Assimilation of GRACE terrestrial water storage
896 data into a land surface model: Results for the Mississippi River basin. *Journal of*



Manuscript in submission to *ESSD*

- 897 *Hydrometeorology*, 9(3), 535-548.
- 898 Zhang, J., Lin, P., Gao, S., & Fang, Z. (2020). Understanding the re-infiltration process to simulating
899 streamflow in North Central Texas using the WRF-hydro modeling system. *Journal of*
900 *Hydrology*, 587, 124902.
- 901 Zhang, J., Wang, T., & Ge, J. (2015). Assessing vegetation cover dynamics induced by policy-driven
902 ecological restoration and implication to soil erosion in southern China. *PLoS One*, 10(6),
903 e0131352.
- 904 Zhang, S., Zhou, L., Zhang, L., Yang, Y., Wei, Z., Zhou, S., Yang, D., Yang, X., Wu, X., & Zhang, Y.
905 (2022). Reconciling disagreement on global river flood changes in a warming climate. *Nature*
906 *Climate Change*, 1-8.
- 907 Zhang, Y., & Liang, S. (2014). Changes in forest biomass and linkage to climate and forest disturbances
908 over Northeastern China. *Global change biology*, 20(8), 2596-2606.
- 909 Zhang, Y., Zheng, H., Zhang, X., Leung, L. R., Liu, C., Zheng, C., Guo, Y., Chiew, F. H., Post, D., &
910 Kong, D. (2023). Future global streamflow declines are probably more severe than previously
911 estimated. *Nature Water*, 1-11.
- 912
- 913