

# CCAM: China Catchment ~~attributes~~Attributes and meteorology for large sample study in contiguous ChinaMeteorology dataset

Zhen Hao<sup>2,\*</sup>, Jin Jin<sup>1,2,\*</sup>, Runliang Xia<sup>2</sup>, Shimin Tian<sup>2</sup>, Wushuang Yang<sup>2</sup>, Qixing Liu<sup>2</sup>, Min Zhu<sup>2</sup>, Tao Ma<sup>2</sup>, Chengran Jing<sup>2</sup>

5 <sup>1</sup>~~School~~<sup>1</sup> School of Computer Science, Northwestern Polytechnical University, Xi'an, China, 710072

<sup>2</sup>~~Yellow~~<sup>2</sup> Yellow River Institute of Hydraulic Research, Zhengzhou, China, 450003

\*These authors contributed equally to this work.

Correspondence to: Jin Jin ([jinjinhao@21cn.com](mailto:jinjinhao@21cn.com))

## Abstract.

10 The lack of a complied large-scale catchment characteristics dataset is a key obstacle limiting the development of large sample hydrology research in China. We introduce the first large-scale catchment attributes ~~and meteorological time series~~-dataset of ~~contiguous~~in China. ~~To develop the dataset, we~~We compiled diverse data sources ~~to generate basin-oriented features describing the catchment characteristics related to hydrological processes. The proposed dataset consists of catchment characteristics, including soil, land cover, climate, topography, and geology, and~~to develop the dataset. The dataset also  
15 includes catchment scale 31-year meteorological time series (from 1990 to 2018). ~~The meteorological variables include precipitation, temperature, evapotranspiration, wind speed, ground surface temperature, pressure, humidity and sunshine duration. We also derived a daily potential~~2020 for each basin. Potential evapotranspiration time series based on ~~a modified Penman's~~Penman's equation- is derived for each basin. The ~~studied~~4911 catchments ~~are~~4875 catchments ~~within contiguous China derived from digital elevation models. We analysed and organised~~included in the spatial variations of catchment  
20 characteristics into a series of maps. Correlation analysis between attributes was conducted. Compared to dataset covers the entire China. We introduced several new indicators describing the catchment geography and the underlying surface compared with previously proposed datasets, we derived more catchment characteristics. The resulting ~~in~~dataset has a total of 125 catchment attributes, providing a complete description of the catchments. Besides, we propose Normal Camels YR, a hydrological. The proposed dataset eoveringalso includes a separate HydroMLYR dataset containing standardized weekly  
25 averaged streamflow for 102 basins ofin the Yellow River ~~basin with normalized streamflow observations.~~Basin. The standardized streamflow data should be able to support machine learning hydrology research in the Yellow River Basin. The proposed dataset ~~provides numerous opportunities for comparative hydrological research, such as examining the difference in hydrological behaviours across different catchments and building general rainfall-runoff modelling frameworks for many catchments instead of limited to a few. The dataset is~~is freely available via <http://doi.org/10.5281/zenodo.4704017> for  
30 ~~community use. We will open source the complement~~at <http://doi.org/10.5281/zenodo.5137288>. In addition, the accompanying code for generating the dataset such that the user can generate meteorological series and catchment attributesis freely available at <https://github.com/haozhen315/CCAM-China-Catchment-Attributes-and-Meteorology-dataset>, supporting

35 the generation of catchment characteristics for any watershed within contiguous China custom basin boundaries. Compiled data for the 4911 basins covering the entire China and the open-sourced code should be able to support the study of any arbitrary basins instead of being limited to only a few basins.

## 1 Introduction

40 Studying a large set of catchments often provides insights that cannot be obtained when looking at a single or few catchments (Coron, Andreassian et al. 2012, Kollat, Reed et al. 2012, Newman, Clark et al. 2015, Lane, Coxon et al. 2019). The hydrologic cycle consists of many sub processes, including evaporation from the ocean, raindrop  
45 evapotranspiration, groundwater flow, subsurface flow and surface runoff, infiltration, etc. are the main components of the terrestrial hydrological cycle. These processes are affected by the nature of the catchment, such as the ability of the soil to hold water. Catchment attributes such as soil characteristics, land cover characteristics and climate indices influence the water movement and storage in these sub processes of the catchment such that hydrologic behaviours can vary across catchments (van Werkhoven, Wagener et al. 2008). The same hydrological model may not be applicable in another basin. However, by examining a large sample of catchments, it is possible for the hydrological model to learn the similarities and differences of hydrological behaviours across catchments. For example, prediction  
50 Studying a large set of terrestrial catchments often provides insights that cannot be obtained when looking at a single or few (Coron, Andreassian et al. 2012, Kollat, Reed et al. 2012, Newman, Clark et al. 2015, Lane, Coxon et al. 2019). For example, a calibrated model may not be applicable in a watershed with vastly different properties. However, by examining a large sample of catchments, it is possible for a data-driven model to learn the similarities and differences of hydrological behaviours across catchments (Kratzert, Klotz et al. 2019). Prediction  
in ungauged basins is a challenging problem present in hydrology. The central challenge is how to extrapolate hydrologic information from gauged basins to ungauged ones. Solving, solving the problem relies on understanding the similarities and differences between different catchments. However, regionally Regionally, and temporally imbalanced observations bring a difficulty to the problem. For a hydrologic model to successfully simulate the ungauged areas, it must adapt itself to the different hydrologic behaviours present in different catchments. (Kratzert, Klotz et al. 2019) Kratzert, Klotz et al. (2019) shows encoding catchment characteristics (e.g., soil characteristics, land cover, topography) into a data-driven model can teachguide the model to behave differently responding to the meteorological time series input based on different sets of static catchment attributes.

60 (Silberstein 2006, Shen, Laloy et al. 2018, Nevo, Anisimov et al. 2019) pointed out that large Large sample hydrological datasets are the foundation and key of many hydrological studies. (Silberstein 2006, Shen, Laloy et al. 2018, Nevo, Anisimov et al. 2019). The term big hydrologic data refers to all data influencing the water cycle, such as the meteorological variables, infiltration characteristics of the study area, land use or land cover types, physical and geological features of the study areacatchment, etc. Many studies cannot be carried out without are based on large-scale hydrologic data (Coron, Andreassian

65 et al. 2012, Singh, van Werkhoven et al. 2014, Berghuijs, Aalbers et al. 2017, Gudmundsson, Leonard et al. 2019, Tyrallis, Papacharalampous et al. 2019). For hydrological research, basin-oriented large-sample datasets are of great significance. For example, comparative hydrology (de Araújo and González Piedra 2009, Singh, Archfield et al. 2014) focus on understanding how hydrological processes interact with the ecosystem, in particular, how hydrologic behaviours change under changes in the surface and sub-surface of the earth to determine to what extent hydrological predictions can be transferred  
70 from one area to another. Large-sample catchment attributes ~~dataset~~datasets provide opportunities for research studying interrelationships among catchment attributes. (~~Seybold, Rothman et al. 2017~~)Seybold, Rothman et al. (2017) studied the correlations between river junction angle with geometric factors, downstream concavity, and aridity. (~~Oudin, Andréassian et al. 2008~~)Oudin, Andréassian et al. (2008) investigates the link between land cover and mean annual streamflow based on 1508 basins representing a large hydroclimatic variety. (~~Voepel, Ruddell et al. 2011~~)Voepel, Ruddell et al. (2011) examines how  
75 the interaction of climate and topography influences vegetation response.

~~Data-driven methods can best benefit from large-scale data. Data-driven approaches have shown great potential in various fields, transforming the applications in many industries (LeCun, Bengio et al. 2015). However, data-driven methods, especially the deep learning-based approaches, usually require high data volumes. Limited data will cause the over-fitting (Blumer, Ehrenfeucht et al. 1987, Abu Mostafa, Magdon Ismail et al. 2012) problem. Therefore, big hydrologic data is the fundamental support for the successful deployment of powerful data-driven strategies.~~

80

~~Traditional hydrological models have some long-standing challenges, such as the inability to capture hydrological processes' mechanism complexity (Kollat, Reed et al. 2012), which is due to the structural limitations of the conceptual models. Data-driven methods are proposed to overcome some existing obstacles. Data-driven strategies open a new way for researchers to acquire knowledge transforming the research pattern from hypothesis-driven to data-driven. (Feng, Fang et al. 2020) proposed a flexible data integration fusing various types of observations to improve rainfall-runoff modelling. The research shows that combining different resources of data benefits predictions in regions with high autocorrelation in streamflow. (Wongso, Nateghi et al. 2020) developed a model predicting the state-level, per capita water uses in the United States, taking various geographic, climatic, and socioeconomic variables as input. The research also identified key factors associated with high water usage. (Mei, Maggioni et al. 2020) proposed a statistical framework for spatial downscaling to obtain hyper-resolution precipitation data. The results show improvements compared with the original product. (Brodeur, Herman et al. 2020) applied machine learning techniques, namely bootstrap aggregation and cross-validation, to reduce overfitting in reservoir control policy search. (Ni and Benson 2020) proposed an unsupervised machine learning method to differentiate flow regimes and identify capillary heterogeneity trapping, showing the promise of machine learning methods for analysing large datasets from coreflooding experiments. (Legasa and Gutiérrez 2020) propose to apply Bayesian Network for multisite precipitation occurrence generation. The proposed methodology shows improvements for existing methods.~~

85

90

95

World-wide data sharing has become a trend (Wickel, Lehner et al. 2007, Ceola, Arheimer et al. 2015, Blume, van Meerveld et al. 2018, Wang, Chen et al. 2020), and the amounts of hydrologic data available are ever-increasing. However, these data typically came from different providers and are compiled in various formats. ~~For example, ASTGTM<sup>+</sup>ASTGTM (Abrams, Crippen et al. 2020) provides a global digital elevation model; Glim (Hartmann and Moosdorf 2012) includes rock types data globally; MODIS provides data products (Knyazikhin 1999, Didan 2015, Myneni, Knyazikhin et al. 2015, Running, Mu et al. 2017, Sulla-Menashe and Friedl 2018) describing features of the land and the atmosphere derived from remote sensing observations; (Yamazaki, Ikeshima et al. 2019) Yamazaki, Ikeshima et al. (2019) provides a global flow direction map at three arc-second resolution; HydroBASINS (Lehner 2014) provides basin boundaries at different scales globally; and GDBD (Masutomi, Inui et al. 2009) provides basin boundaries with geographic attributes; GLHYMPS (Gleeson, Moosdorf et al. 2014) provides a global map of subsurface permeability and porosity; SoilGrids250m (Hengl, Mendes de Jesus et al. 2017) dataset provides global numeric soil properties. Local government agencies often hold meteorological data such as precipitation and evaporation, and the amount of this data is also growing, however, data transparency has still been a problem (Viglione, Borga et al. 2010). The~~

~~However, the~~ data mentioned above are rarely spatially aggregated to the catchment-scale, making it difficult for researchers to use these data. Properly pre-processed and formatted datasets ~~on a large scale~~ are of great importance for ~~the~~ hydrology research. Searching for appropriate data sources, pre-processing, and formatting often consumes a lot of ~~researchers'~~ time. In some cases, individual research groups either do not know where to obtain the appropriate data or cannot properly process the data to receive the desired format. ~~In summary, although data sharing is being advocated in the community, it is usually difficult for the public to obtain the required data, either because there are not enough observations or because of the difficulties in the data processing.~~

~~In summary, both data-driven and traditional hydrological research need diverse hydrologic datasets to learn the generalisation capability from one area to another. For a model to adapt to various behaviours in different catchments, the dataset must be large enough to represent the complex heterogeneity presented in the natural hydrologic system. Although data sharing is being advocated in the community, it is usually difficult for the public to obtain certain data such as meteorological data and streamflow observations, either because there are not enough observations or because there are no open access permissions.~~

Recently, there are efforts (Addor, Newman et al. 2017, Alvarez-Garreton, Mendoza et al. 2018, Chagas, Chaffe et al. 2020, Coxon, Addor et al. 2020) ~~compiling different types of data sources to form large scale hydrological datasets. These four collected datasets cover the continental United States, Chile, Brazil, and Great Britain. (Addor, Do et al. 2020) reviewed these datasets and discussed the guidelines for producing large sample hydrological datasets and the limitations of the currently proposed datasets. The CAMELS dataset has been used to support a lot of research. Based on CAMELS, (Kratzert, Klotz et~~

---

<sup>+</sup> <https://asterweb.jpl.nasa.gov/gdem.asp>



al. 2018) built a Long Short-Term Memory (LSTM) network for rainfall-runoff modelling, showing that one model can predict the discharge for a variety of catchments. (Knoben, Freer et al. 2019) compared metrics used in hydrology based on simulations on many basins. (Tyrallis, Papacharalampous et al. 2019) studied the relationship between the shape parameter and basin attributes based on the sizeable basin-oriented dataset.

135

However, there is no large-scale compilation of hydrological datasets in contiguous China. An alternative is on a global scale, the HydroATLAS (Linke, Lehner et al. 2019) dataset. However, since it is on a world-wide scale, compared with other datasets constructed for regions, the dataset lacks many attributes and is not built according to the CAMELS standards. Besides, the climatic data is not up to date (1950-2000), and the derivation of climatic data lacks ground surface observations inputs, such

140

that the data quality is not guaranteed.

Therefore, researchers still need to do repetitive works to compile data from different sources such as obtaining historical meteorological data (temperature, rainfall, evapotranspiration) of a catchment in contiguous China. Inspired by (Addor, Newman et al. 2017), in this paper, we present a catchment scale hydrologic dataset compiling a wide variety of hydrological data, including basin topography, climate indices, land cover characteristics, soil characteristics and geological characteristics covering contiguous China.

145

to compile different types of data sources forming large scale hydrological datasets. These four collected datasets cover the continental United States, Chile, Brazil, and Great Britain. Addor, Do et al. (2020) reviewed these datasets and discussed the guidelines for producing large-sample hydrological datasets and the limitations of the currently proposed datasets. The static properties of 671 river basins in the United States are calculated by CAMELS (Addor, Newman et al. 2017), which is an extension of a previously proposed hydrometeorological data set (Newman, Clark et al. 2015). Unfortunately, it is impossible to publish streamflow data in China for the time being. The CAMELS dataset has been used to support a lot of research. For example, Knoben, Freer et al. (2019) compared metrics used in hydrology based on simulations on many basins. Tyrallis, Papacharalampous et al. (2019) studied the relationship between the shape parameter and basin attributes based on the sizeable basin-oriented dataset.

150

155

There is currently no compilation of China-specific catchment attributes datasets. An alternative, the HydroATLAS (Linke, Lehner et al. 2019) dataset, which is on a global scale, is basically performing zonal statistics on the source data. HydroATLAS lacks many indicators which need derivations based on the source data, such as rainfall seasonality, the fraction of precipitation falling as snow, basin shape factors and root depth distributions. What's worse, the meteorological data is only up to 2000, which is outdated.

160

In summary, a lack of a compiled catchment attributes dataset is a key obstacle limiting the development of large sample hydrology research in China. Inspired by (Addor, Newman et al. 2017), we compiled multiple data sources, including basin topography, climate indices, land cover characteristics, soil characteristics and geological characteristics. Different from

165 (Addor, Newman et al. 2017), the catchments included in the dataset covers the entire study area, instead of being limited to a few.

The proposed dataset is the first dataset providing catchments meteorological time series and catchments attributes of ~~contiguous~~ China. We compiled and named the dataset following most standards of the previously proposed datasets. ~~Unlike CAMELS and CAMELS-CL, catchments in the proposed dataset are not selective. Instead, the~~The dataset consists of all ~~generated basins derived basin boundaries~~ from the Digital Elevation Model (DEM), ~~based on~~ which came from the Global Drainage Basin Dataset (Masutomi, Inui et al. 2009). The Global Drainage Basin Dataset (GDBD) is derived at high-resolution (100m-1km) and has a good geographic agreement with existing global drainage basin data in China<sup>2</sup>. ~~Besides, an essential feature of the. In addition, previously proposed dataset is that it provides a complete description of the catchment, rather than an abstraction. For example, both CAMELS and CAMELS-CL only datasets (Addor, Newman et al. 2017, Alvarez-Garreton,~~  
170 ~~Mendoza et al. 2018, Chagas, Chaffe et al. 2020, Coxon, Addor et al. 2020) report only the most frequent and second most frequent catchment land cover and lithology types. Instead, the proposed dataset~~CCAM calculates the ~~proportion~~proportions of ~~each~~all land cover and lithology ~~type for each catchment to serve data-driven research better. We also introduced many more climate characteristics and soil characteristics to support more diverse potential research~~types.

180 ~~Researchers from different places can use the proposed dataset in conjunction with their streamflow data, simplifying organising and compiling various data resources, which is usually repetitive work. The proposed dataset is undoubtedly the most comprehensive catchment attributes and meteorological time series dataset in contiguous China and is suitable for multi-purpose data-driven research. The dataset consists of basin boundaries in the shapefile format, computed catchment attributes of climate, land cover, soil, topography and lithology and 29 year meteorological time series. Table 1 compares the number of static attributes between CAMELS, CAMELS-BR, and the proposed dataset.~~

In addition to the basin-wise attributes provided in CCAM, we propose HydroMLYR, a hydrology dataset for machine learning research in the Yellow River Basin providing weekly averaged standardized streamflow data for 102 basins in the Yellow River Basin (YRB). HydroMLYR is proposed to support machine learning hydrology research at YRB. Traditional  
190 hydrological models have some long standing challenges, such as the inability to capture hydrological processes' mechanism complexity (Kollat, Reed et al. 2012), which is due to the structural limitations of the conceptual models. Data-driven strategies represented by machine learning are proposed to overcome some existing obstacles and they open a new way for researchers to acquire knowledge transforming the research pattern from hypothesis-driven to data-driven. Feng, Fang et al. (2020)

---

<sup>2</sup> ~~In this study, gauge streamflow measurements are not available in areas other than the Yellow River such that it is infeasible to specify a gauge location for generating the basin boundary for most of the areas. Streamflow measurements have strict redistribution policy; however, local research institutions have their streamflow measurements for hydrological research, the proposed dataset can used in conjunction with the streamflow data of researchers in various places.~~

195 proposed a flexible data integration fusing various types of observations to improve rainfall-runoff modelling. The research shows that combining different resources of data benefits predictions in regions with high autocorrelation in streamflow. Wongso, Nateghi et al. (2020) developed a model predicting the state-level, per capita water uses in the United States, taking various geographic, climatic, and socioeconomic variables as input. The research also identified key factors associated with high water usage. Mei, Maggioni et al. (2020) proposed a statistical framework for spatial downscaling to obtain hyper-resolution precipitation data. The results show improvements compared with the original product. Brodeur, Herman et al. (2020) applied machine learning techniques, namely bootstrap aggregation and cross-validation, to reduce overfitting in reservoir control policy search. Ni and Benson (2020) proposed an unsupervised machine learning method to differentiate flow regimes and identify capillary heterogeneity trapping, showing the promise of machine learning methods for analysing large datasets from coreflooding experiments. Legasa and Gutiérrez (2020) propose to apply Bayesian Network for multisite precipitation occurrence generation, and the proposed methodology shows improvements for existing methods. The proposed data set can be used to develop or verify machine learning models in the YRB.

200 The paper is organized as follows: Section 2 describes the study area. Section 3-7 describes the five classes of the computed catchment attributes. ~~In section 3-7, each unit follows the same structure: first introduce the meaning and significance of each added feature and data source used, then describe the variables' spatial variability if necessary.~~ Section 8 describes the proposed catchment-scale meteorological ~~forcing~~-time series. Section 9 introduce the ~~Normal Camels YRHydroMLYR~~ dataset, ~~which provides normalized streamflow measurements for 102 catchments of Yellow River.~~ Section 10 describes the code and data availability. Section 11 ~~presents~~is the concluding remark.

In summary, our contributions are as follows:

- 215 (1) ~~The proposed dataset is the first large scale dataset containing catchment scale meteorological time series of contiguous China, which is the basis for many hydrological studies.~~
- (2) ~~We present the first basin-oriented static attributes dataset in contiguous China.~~
- (3) ~~We introduce several new catchment characteristics providing a complete description of the catchment compared with the previously proposed datasets such that the proposed dataset is prepared for potential hydrological studies.~~
- 220 (4) ~~We offer a self-contained dataset covering 102 basins of the Yellow River basin with normalized runoff observation supporting many potential studies.~~
- (5) ~~We will open source the code for generating the dataset such that the user can generate a dataset for any watershed within contiguous China.~~

**Table 1 Number of computed attributes in CAMELS, CAMELS-BR and the proposed dataset.**

Attribute class	CAMELS(A17)	CAMELS-BR	Ours
<del>Location and topography</del>	<del>9</del>	<del>11</del>	<del>12</del>

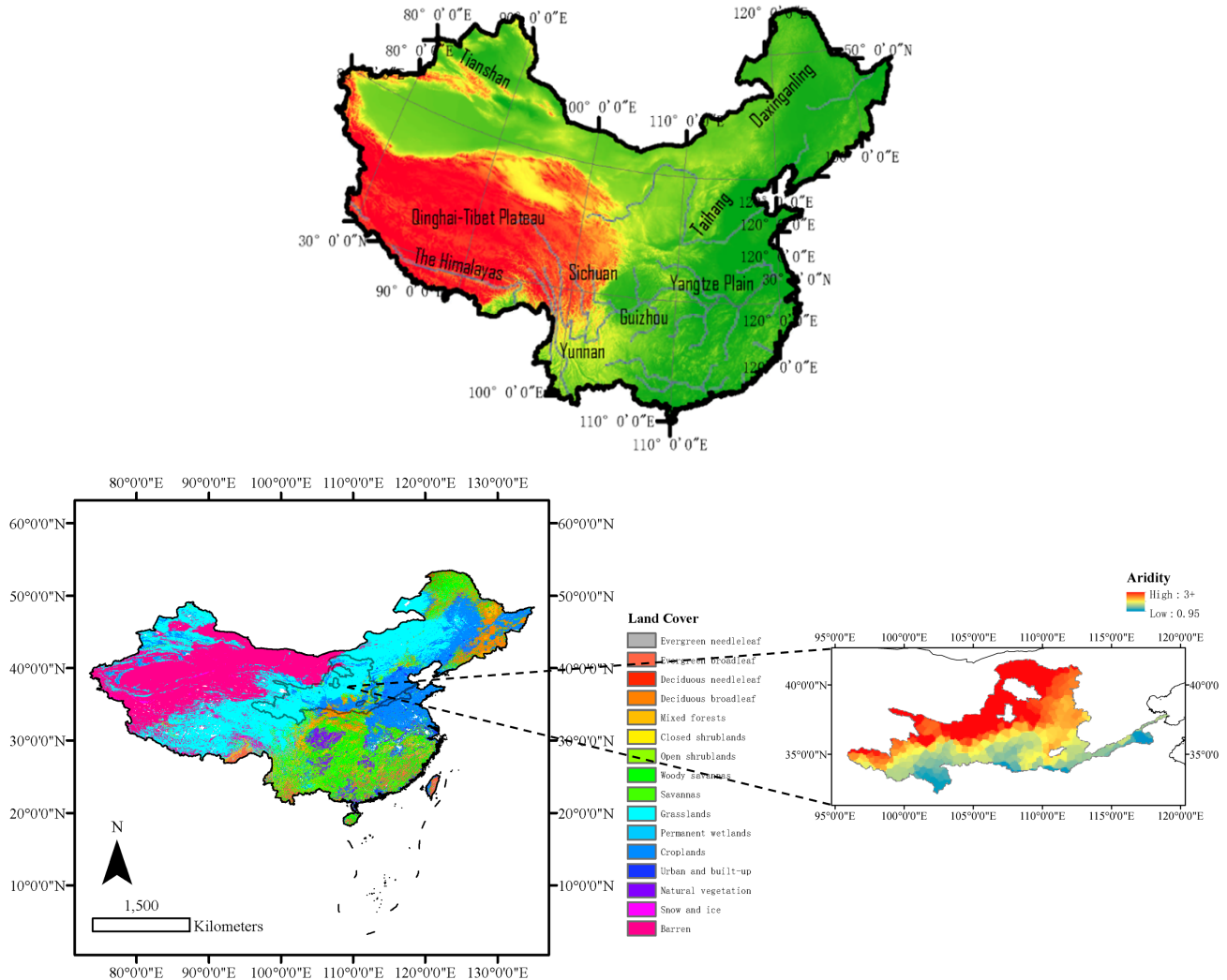
<del>Geology</del>	<del>7</del>	<del>7</del>	<del>18</del>
<del>Soil</del>	<del>11</del>	<del>6</del>	<del>54</del>
<del>Land cover</del>	<del>8</del>	<del>11</del>	<del>22</del>
<del>Climatic indices</del>	<del>11</del>	<del>13</del>	<del>17</del>
<del>Human intervention indices</del>	<del>not computed</del>	<del>4</del>	<del>2</del>
<del>Total</del>	<del>46</del>	<del>52</del>	<del>125</del>

225

**Table 2 Summary of basin daily discharge and forcing data in CAMELS, CAMELS-BR and the proposed dataset.**

Forcing data class	CAMELS	CAMELS-BR	Ours
Temperature	available	available	available
Precipitation	available	available	available
Solar radiation	available	not available	available
Day length	available	not available	not available
Sunshine hours	not available	not available	available
Humidity	available	not available	available
Snow water equivalent	available	not available	not available
Wind velocity	not available	not available	available
Ground surface pressure	available	not available	available
Observed evaporation	not available	available	available
Potential evapotranspiration	not available	available	available
Streamflow	available	available	partially available (see Section 9)

## 2 Study area



(a)

(

b

)

**Figure 1.** Overview of the study area. The study area covers a wide range of latitude and longitude, from 18.2° N to 52.3° N, and from 76.0° E to 134.3° E. (a) The main geographical features map of contiguous China. China is mountainous; mountains and hills occupy two-thirds of the area. (b) The distribution map of the delimited catchments based on the ASTER DEM, the catchments studied are all catchment areas delimited from the DEM, covering contiguous China, with 4875 catchments, most of which are 2000 to 5000 square kilometres.

: Left: Study area of CCAM and the distribution of land cover types. The studied basins cover the whole of China. Right: Study area of HydroMLYR and the distribution of aridity (PET/P) index. YRB is a generally arid area. The data set provided can be used as a good sample for studying hydrology in arid regions.

The study area corresponds to ~~contiguous~~the whole of China; (Fig. 1), with diverse climate and terrain characteristics, spanning from 18.2° N to 52.3° N and 76.0° E to 134.3° E. Mountains, plateaus, and hills account for about two-thirds of areas of ~~contiguous~~China, and the remaining are basins and plains. China's topography is like a three-level ladder, high in the west and low in the east. The Qinghai-Tibet Plateau, the highest plateau globally, located in the west of ~~contiguous~~China, with a mean elevation of over 4000 meters, is the first step of China's topography. The Xinjiang region, the Loess Plateau, the Sichuan Basin, and the Yunnan-Guizhou Plateau to the north and east are the second step of China's topography. The mean sea level here is between 1000 to 2000 meters. Plains and hills dominate the east of the Daxinganling-Taihang Mountain to the coastline, the third step of ~~contiguous~~China. The elevation of this step descends to 500-1,000 meters. To better characterize the studied catchments, we have derived various attributes. Table 1 compares the number of derived attributes between several proposed datasets.

**Table 1: Number of computed attributes in CAMELS, CAMELS-BR and CCAM.**

<u>Attribute class</u>	<u>CAMELS(A17)</u>	<u>CAMELS-BR</u>	<u>CCAM</u>
<u>Location and topography</u>	<u>9</u>	<u>11</u>	<u>12</u>
<u>Geology</u>	<u>7</u>	<u>7</u>	<u>18</u>
<u>Soil</u>	<u>11</u>	<u>6</u>	<u>54</u>
<u>Land cover</u>	<u>8</u>	<u>11</u>	<u>22</u>
<u>Climatic indices</u>	<u>11</u>	<u>13</u>	<u>17</u>
<u>Human intervention indices</u>	<u>not computed</u>	<u>4</u>	<u>2</u>
<u>Total</u>	<u>46</u>	<u>52</u>	<u>125</u>

In ~~contiguous~~ China, precipitation and temperature vary significantly in different places, forming a diverse climate environment. According to the Köppen Climate Classification System, from northwest to southeast, China's climate gradually evolves from Cold desert (BW<sub>k</sub>) climate, Tundra (ET) climate, Warm and temperate continental (D<sub>fa</sub> and D<sub>wb</sub>) climate to Humid subtropical (C<sub>wa</sub>) climate and Warm oceanic (C<sub>fa</sub>) climate. From the perspective of temperature zones, there are tropical, subtropical, warm temperate, medium temperate, cold temperate and Qinghai-Tibet Plateau regions, and there are humid regions, semi-humid regions, semiarid regions, and arid regions from the perspective of wet and dry zones. Moreover, the same temperature zone can contain different dry and wet zones. Therefore, there will be differences in heat and wetness in the same climate type. The complexity of the terrain makes the climate even more complex and diverse. Besides, China has a wide range of regions affected by the alternating winter and summer monsoons. Compared with other parts of the world at the same latitude, these areas have low winter temperatures, high summer temperatures, significant annual temperature differences, and concentrated precipitation in summer. The cold and dry winter monsoon occurs in Asia's interior, far away from the ocean. Under its influence, winter rainfall in most parts of China is low, accompanied by low temperature. The summer monsoon is

warm and humid, coming from the Pacific Ocean and the Indian Ocean. Under its influence, precipitation generally increases.

Table 2 compares the provided forcing variables in CAMELS, CAMELS-BR and CCAM.

**Table 2: Summary of forcing variables provided in CAMELS, CAMELS-BR and CCAM.**

<u>Forcing data class</u>	<u>CAMELS</u>	<u>CAMELS-BR</u>	<u>CCAM</u>
<u>Temperature</u>	<u>available</u>	<u>available</u>	<u>available</u>
<u>Precipitation</u>	<u>available</u>	<u>available</u>	<u>available</u>
<u>Solar radiation</u>	<u>available</u>	<u>not available</u>	<u>available</u>
<u>Day length</u>	<u>available</u>	<u>not available</u>	<u>not available</u>
<u>Sunshine hours</u>	<u>not available</u>	<u>not available</u>	<u>available</u>
<u>Humidity</u>	<u>available</u>	<u>not available</u>	<u>available</u>
<u>Snow water equivalent</u>	<u>available</u>	<u>not available</u>	<u>not available</u>
<u>Wind velocity</u>	<u>not available</u>	<u>not available</u>	<u>available</u>
<u>Ground surface pressure</u>	<u>available</u>	<u>not available</u>	<u>available</u>
<u>Observed evaporation</u>	<u>not available</u>	<u>available</u>	<u>available</u>
<u>Potential evapotranspiration</u>	<u>not available</u>	<u>available</u>	<u>available</u>

265

**Table 3: Summary table of catchment attributes available in the proposed dataset.**

<u>Attribute class</u>	<u>Attribute name</u>	<u>Description</u>	<u>Unit</u>	<u>Data source</u>
Climate indices (computed for 1 Oct 1990 to 30 Sep 2018)	pet_mean	mean daily pet (Penman–Monteith equation)	mm d <sup>-1</sup>	(Subramanya 2013)
	evp_mean	mean daily evaporation (observations)	mm d <sup>-1</sup>	SURF_CLI_CHN_MUL_DAY3F <sup>3</sup>
	gst_mean	mean daily ground surface temperature	°C	
	pre_mean	mean daily precipitation	mm d <sup>-1</sup>	
	prs_mean	mean daily ground surface pressure	hPa	
	rhu_mean	mean daily relative humidity	-	
	ssd_mean	mean daily sunshine duration	h	

<sup>3</sup> [http://data.cma.cn/data/cdcdetail/dataCode/SURF\\_CLI\\_CHN\\_MUL\\_DAY.html](http://data.cma.cn/data/cdcdetail/dataCode/SURF_CLI_CHN_MUL_DAY.html)



tem_mean	mean daily temperature	°C
win_mean	mean daily wind speed	m s <sup>-1</sup>
p_seasonality	seasonality and timing of precipitation (estimated using sine curves to represent the annual temperature and precipitation cycles, positive [negative] values indicate that precipitation peaks in summer [winter], values close to 0 indicate uniform precipitation throughout the year)	-
high_prec_freq	frequency of high-precipitation days ( $\geq 5$ times mean daily precipitation)	d yr <sup>-1</sup>
high_prec_dur	average duration of high-precipitation events (number of consecutive days $\geq 5$ times mean daily precipitation)	d
high_prec_timing	season during which most high-precipitation days ( $\geq 5$ times mean daily precipitation) occur	season
low_prec_freq	frequency of dry days ( $< 1$ mm d <sup>-1</sup> )	d yr <sup>-1</sup>
low_prec_dur	average duration of dry periods (number of consecutive days $< 1$ mm d <sup>-1</sup> )	d
low_prec_timing	season during which most dry days ( $< 1$ mm d <sup>-1</sup> ) occur	season
frac_snow_daily	fraction of precipitation falling as snow (for days colder than 0 °C)	-
p_seasonality	seasonality and timing of precipitation, positive [negative] values indicate that precipitation peaks in summer [winter], values	-

		close to 0 indicate uniform precipitation throughout the year		
Geological characteristics	geol_porosity	subsurface porosity	-	(Gleeson, Moosdorf et al. 2014)
	geol_permeability	subsurface permeability (log-10)	m <sup>2</sup>	
	ig	fraction of the catchment area associated with ice and glaciers	-	(Hartmann and Moosdorf 2012)
	pa	fraction of the catchment area associated with acid plutonic rocks	-	
	sc	fraction of the catchment area associated with carbonate sedimentary rocks	-	
	su	fraction of the catchment area associated with unconsolidated sediments	-	
	sm	fraction of the catchment area associated with mixed sedimentary rocks	-	
	vi	fraction of the catchment area associated with intermediate volcanic rocks	-	
	mt	fraction of the catchment area associated with metamorphic	-	
	ss	fraction of the catchment area associated with siliciclastic sedimentary rocks	-	
pi	fraction of the catchment area associated with intermediate plutonic rocks	-		
va	fraction of the catchment area associated with acid volcanic rocks	-		
wb	fraction of the catchment area associated with water bodies	-		

	pb	fraction of the catchment area - associated with basic plutonic rocks	
	vb	fraction of the catchment area - associated with basic volcanic rocks	
	nd	fraction of the catchment area - associated with no data	
	py	fraction of the catchment area - associated with pyroclastic	
	ev	fraction of the catchment area - associated with evaporites	
Land cover characteristics	lai_max	maximum monthly mean of the leaf area index (based on 12 monthly means)	(Myneni, Knyazikhin et al. 2015)
	lai_diff	difference between the maximum and minimum monthly mean of the leaf area index (based on 12 monthly means)	
	ndvi_mean	mean normalized difference vegetation index (NDVI)	(Didan 2015)
	root_depth_50	root depth (percentiles=50% extracted from a root depth distribution based on IGBP land cover)	m Eq. 2 and Table 2 in (Zeng 2001)
	root_depth_99	root depth (percentiles=99% extracted from a root depth distribution based on IGBP land cover)	m
	evergreen needleleaf tree	catchment area fraction covered by evergreen needleleaf tree	-
evergreen broadleaf tree	catchment area fraction covered by evergreen broadleaf tree	-	

	deciduous needleleaf tree	catchment area fraction covered - by deciduous needleleaf forests		
	deciduous broadleaf tree	catchment area fraction covered - by deciduous broadleaf tree		
	mixed forest	catchment area fraction covered - by mixed forest		
	closed shrubland	catchment area fraction covered - by closed shrubland		
	open shrubland	catchment area fraction covered - by open shrubland		
	woody savanna	catchment area fraction covered - by woody savanna		
	savanna	catchment area fraction covered - by savanna		
	grassland	catchment area fraction covered - by grassland		
	permanent wetland	catchment area fraction covered - by permanent wetland		
	cropland	catchment area fraction covered - by cropland		
	urban and built-up land	catchment area fraction covered - by urban and built-up land		
	cropland/natural vegetation	catchment area fraction covered - by cropland/natural vegetation		
	snow and ice	catchment area fraction covered - by snow and ice		
	barren	catchment area fraction covered - by barren		
	water bodies	catchment area fraction covered - by water bodies		
Topography, location, and	basin_id pop	drainage basin identifiers population	- people	(Masutomi, Inui et al. 2009)

Human			people	
intervention	pop_dnsty	population density	km <sup>-2</sup>	
	lat	mean latitude	°N	
	lon	mean longitude	°E	
	elev	mean elevation	M	
	area	catchment area	km <sup>2</sup>	
				m km <sup>-1</sup> (Horn 1981)
	slope	mean slope		1
	length	The length of the mainstream measured from the basin outlet to the remotest point on the basin boundary. The mainstream is identified by starting from the basin outlet and moving up the catchment.	km	km (Subramanya 2013)
	form factor	catchment area / (catchment length) <sup>2</sup>		-
	shape factor	(catchment length) <sup>2</sup> / catchment area		-
compactness coefficient	perimeter of the catchment / perimeter of the circle whose area is that of the basin		-	
circulatory ratio	catchment area / area of circle of catchment perimeter		-	
elongation ratio	diameter of circle whose area is basin area / catchment length		-	
Soil	pdep	soil profile depth	cm	(Shangguan, Dai et al. 2013)
	clay	percentage of clay content of the soil material	%	
	sand	percentage of sand content of the soil material	%	
	por	porosity	cm <sup>3</sup> cm <sup>-3</sup>	

silt	percentage of silt content of the soil material	%	
grav	rock fragment content	%	
som	soil organic carbon content	%	
log_k_s4F <sup>4</sup>	log-10 transformation of saturated hydraulic conductivity	cm d <sup>-1</sup>	(Dai, Xin et al. 2019)
theta_s <sup>4</sup>	saturated water content	cm <sup>3</sup> cm <sup>-3</sup>	
tkstatu <sup>4</sup>	thermal conductivity of unfrozen saturated soils	W m <sup>-1</sup> K <sup>-1</sup>	
bldfie <sup>4</sup>	bulk density	kg m <sup>-3</sup>	(Hengl, Mendes de Jesus et al. 2017)
cecsol <sup>4</sup>	cation-exchange capacity	cmol+ kg <sup>-1</sup>	
orcdrc <sup>4</sup>	organic carbon content	g kg <sup>-1</sup>	
phihox <sup>4</sup>	pH in H2O	10 <sup>-1</sup>	
bdticm	depth to bedrock	cm	

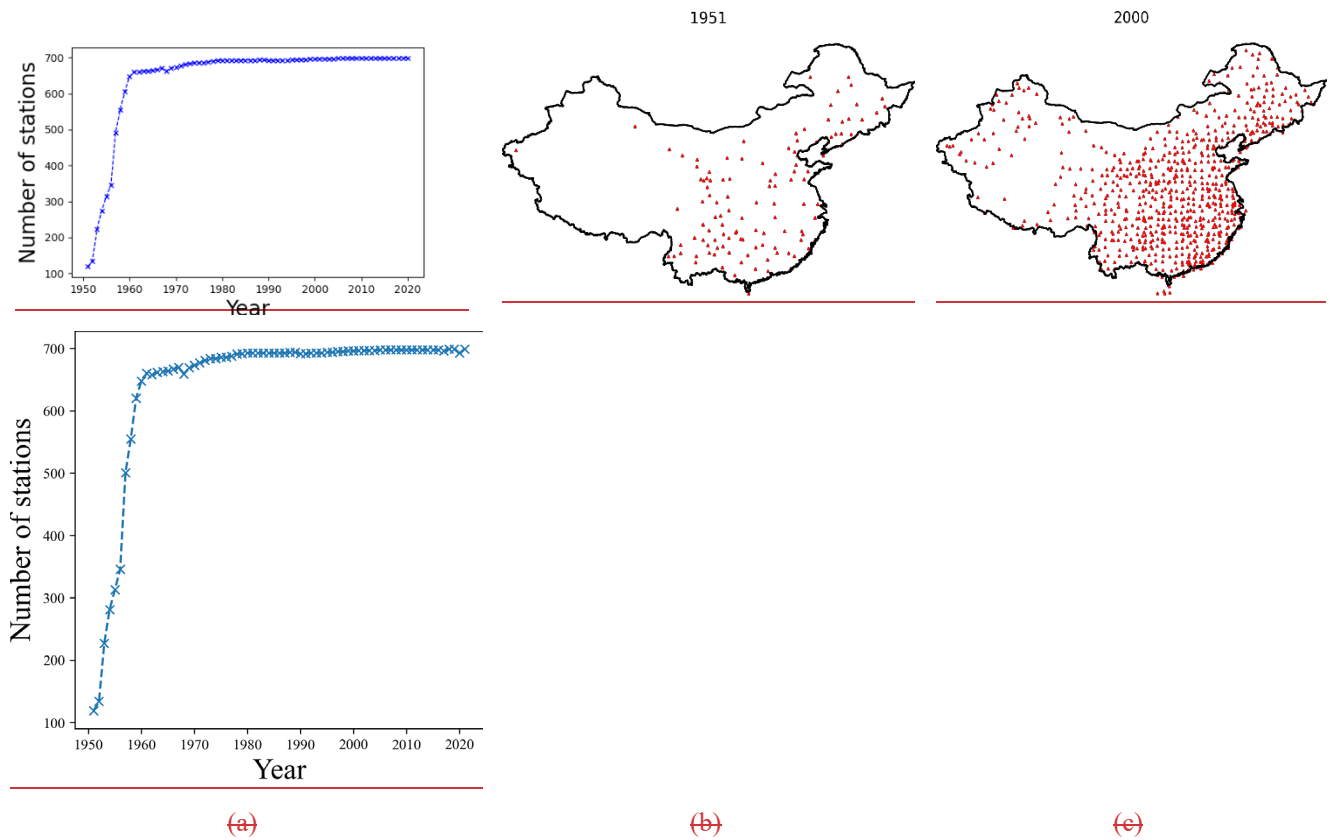
### 3 ~~Climate~~Climatic indices

270 ~~Meteorological raw~~Raw meteorological data ~~was~~is provided by the China Meteorological Data ~~Network~~<sup>3</sup>Network, released as the SURF\_CLI\_CHN\_MUL\_DAY (V3.0) dataset<sup>5</sup>, which provides ~~complete variable types and~~the longest period (1951-~~2018~~2020) of meteorological time series ~~of~~in China. The SURF\_CLI\_CHN\_MUL\_DAY product includes site observations of pressure, temperature, relative humidity, precipitation, evaporation, wind speed, sunshine duration, and ground surface temperature. ~~The summary is presented in Table 4. (Table 4).~~ The Inverse distance weighting method is used for interpolating the site observations. ~~Climate indices are then obtained by taking the average of the catchment scale extraction from the interpolated raster.~~ To ensure data quality, we ~~choose~~use the latter ~~2931~~-year record (from 1990 to ~~2018~~2020) to construct the dataset since sites' distribution was sparse in the early days (Fig. 2). We computed more climatic characteristics compared with other datasets (Table 2). These ~~characteristics have critical potential effects on the~~variables are useful in hydrological ~~processes~~modelling; for example, wind speed can affect actual evapotranspiration. To be consistent with the CAMELS (Addor, Newman et al. 2017), we ~~also~~ determined all climatic attributes (Woods 2009) ~~provided~~in the CAMELS dataset. ~~The~~As a

<sup>4</sup> The data source contains multi-layer soil data, soil characteristics for all layers are determined.

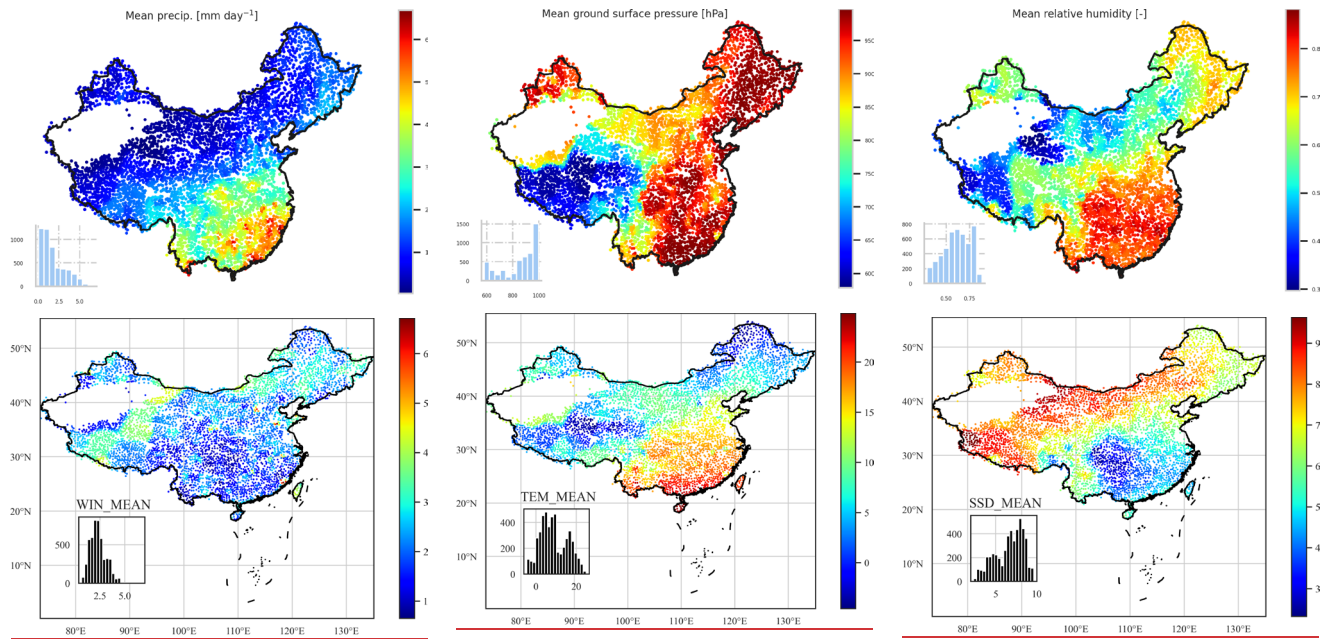
<sup>5</sup> SURF\_CLI\_CHN\_MUL\_DAY is freely available for global researchers.

280 result, the proposed dataset provides more meteorological variables and longer time series (1990-~~2018~~2020) than CAMELS and CAMELS-CL. A summary of the ~~computed Climate~~derived climate indices is presented in Table 3. The national ~~distribution~~distributions of ~~meteorological attributes of catchments is~~the climate indicators are shown in ~~Fig. 3.~~Fig. 3.



285 **Figure 2. Overview of changes: Changes** in the number ~~and distribution~~ of meteorological stations in China. ~~(a) The number of meteorological stations varies with the year.~~ There were only 119 stations in 1951. This number increased rapidly from 1951 to the early 1960s, and the number of stations remained stable after 2000. ~~(b) Distribution map of China's meteorological stations in 1951.~~ ~~(c) Distribution map of China's meteorological stations in 2000.~~ To ensure the data quality, we used the latter 31-year records (from 1990 to 2020) to construct the dataset.

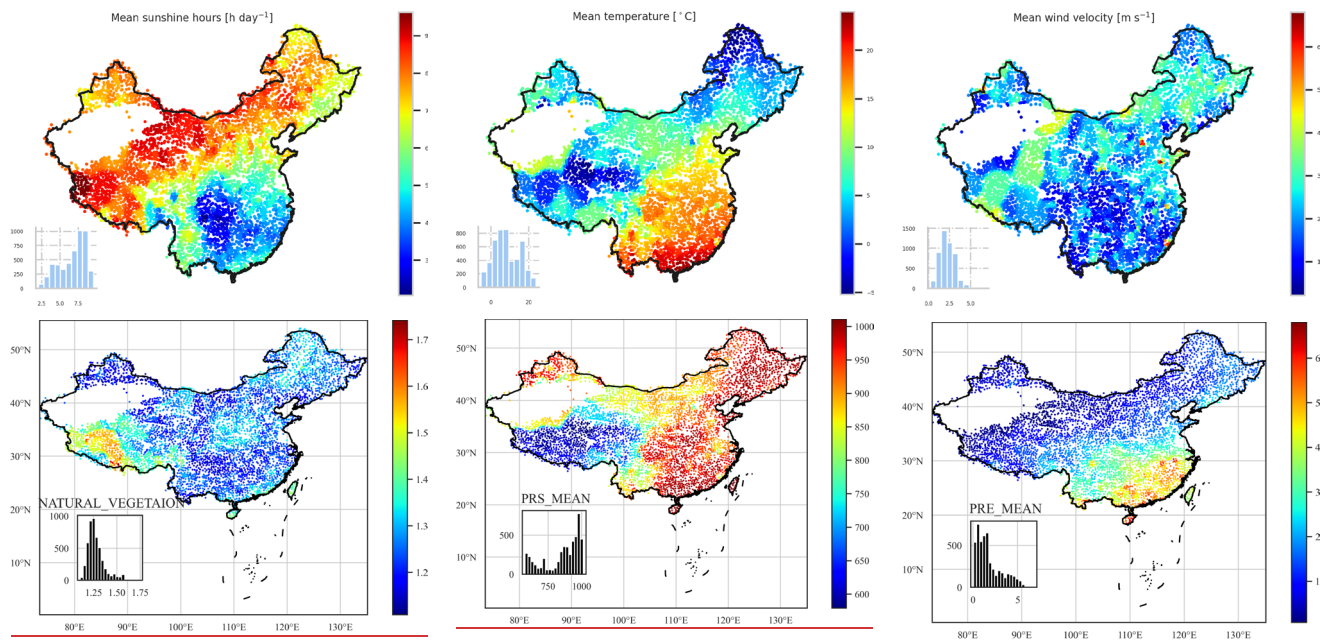




(a)

(b)

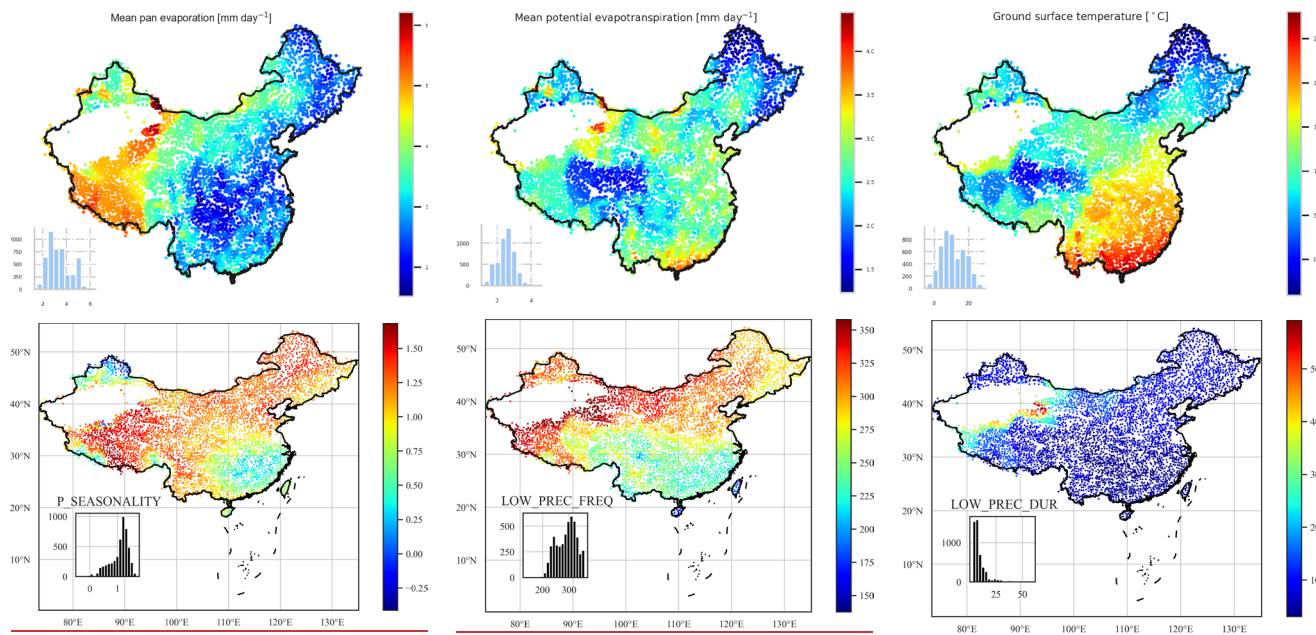
(c)



(d)

(e)

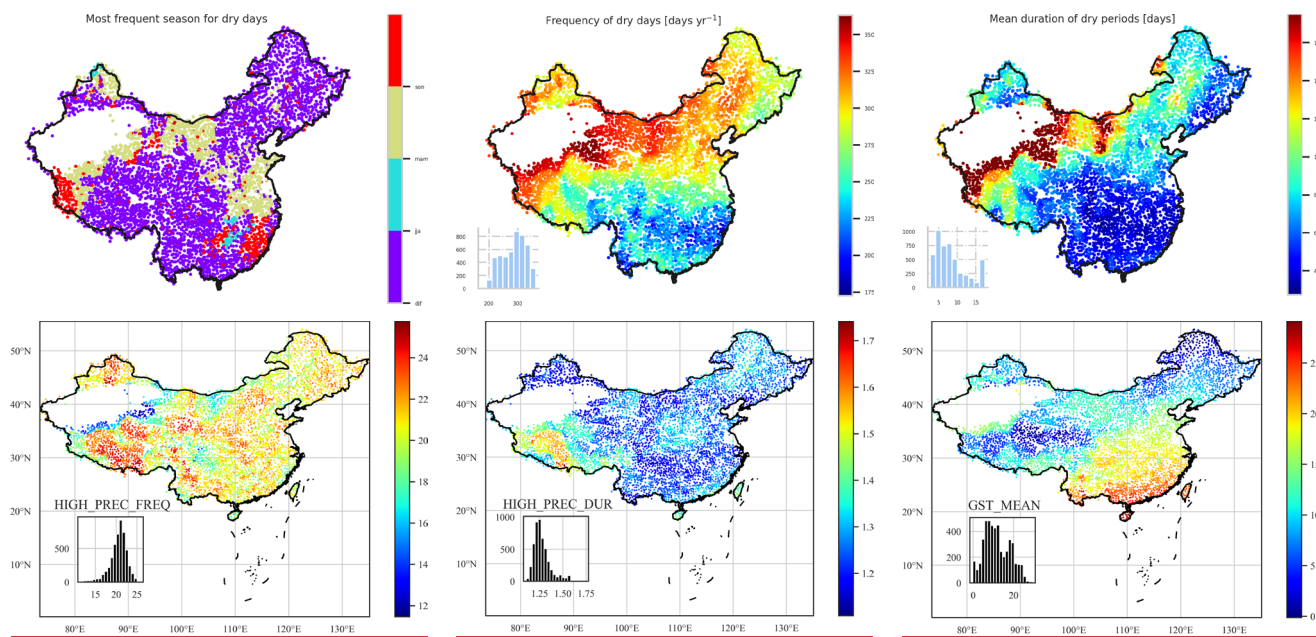
(f)



(g)

(h)

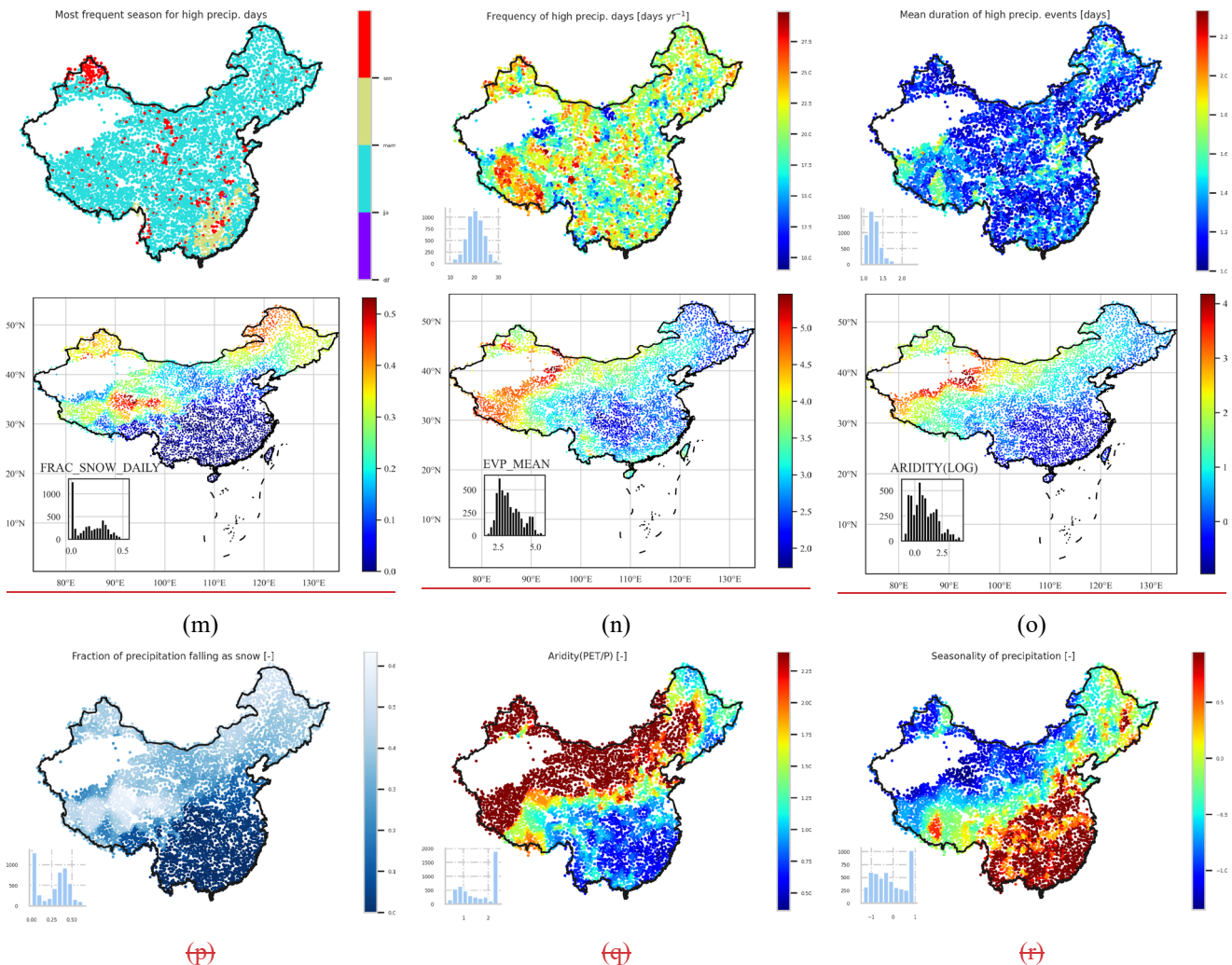
(i)



(j)

(k)

(l)



**Figure 3-Maps: Distributions of climatic indices over contiguous China. The histograms and bar-plots indicate All basins are plotted in the numbersame size. When extreme values of catchments (out of 4875) in each bin or categorya variable affect visualization (cause most areas to have the same colour), the log values are used for visualization.**

290

The instruments for measuring potential evaporation were updated from 2000 to 2005. Early observations can be multiplied by a correction coefficient to approximate the new tools. However, the coefficient varies across stations making the approach infeasible. To complement this, we calculated potential evapotranspiration (PET) based on a modified Penman's Equation (see Appendix A) and other observed meteorological variables, providing a series of consistent evapotranspiration estimationpotential evaporation estimations for reference.

295

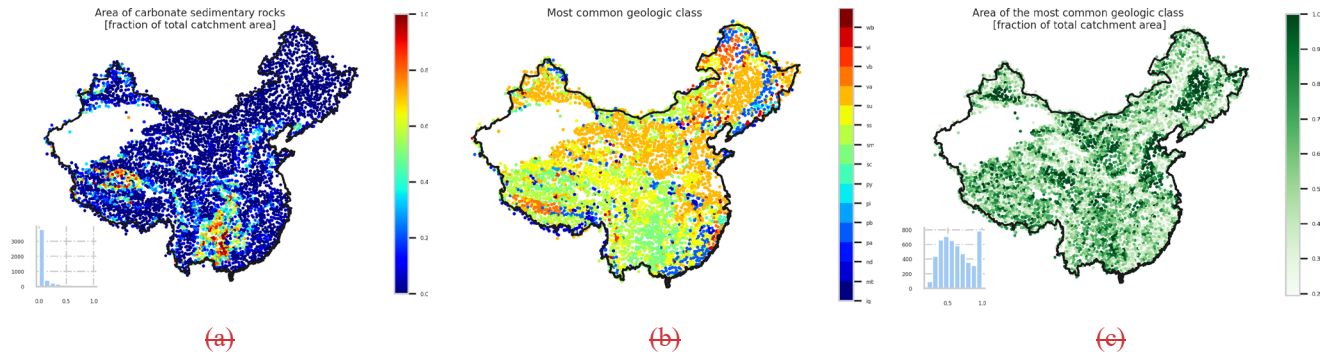
The average daily precipitation in contiguous China is highest in the southeast and lowest in the northwest. It is also higher in the coastal areas than in the interior land. Ground surface pressure is positively correlated with elevation, the highest in the Qinghai-Tibet Plateau and the lowest in the Southeast Plain. The average relative humidity is generally positively correlated

300 with precipitation; they are also higher in some forested areas, such as the Taihang Mountains and Daxingan Mountains. The Qinghai-Tibet Plateau has the lowest average temperature, and the southern coastal area has the highest. A distinctive feature of the distribution of wind speed is the high wind speed in mountainous areas. The highest wind speed occurs in the southeast coastal area (> 6 meters per second). ~~Refer to Section 8 for a detailed description of the proposed catchment-scale meteorological time series dataset of contiguous China.~~

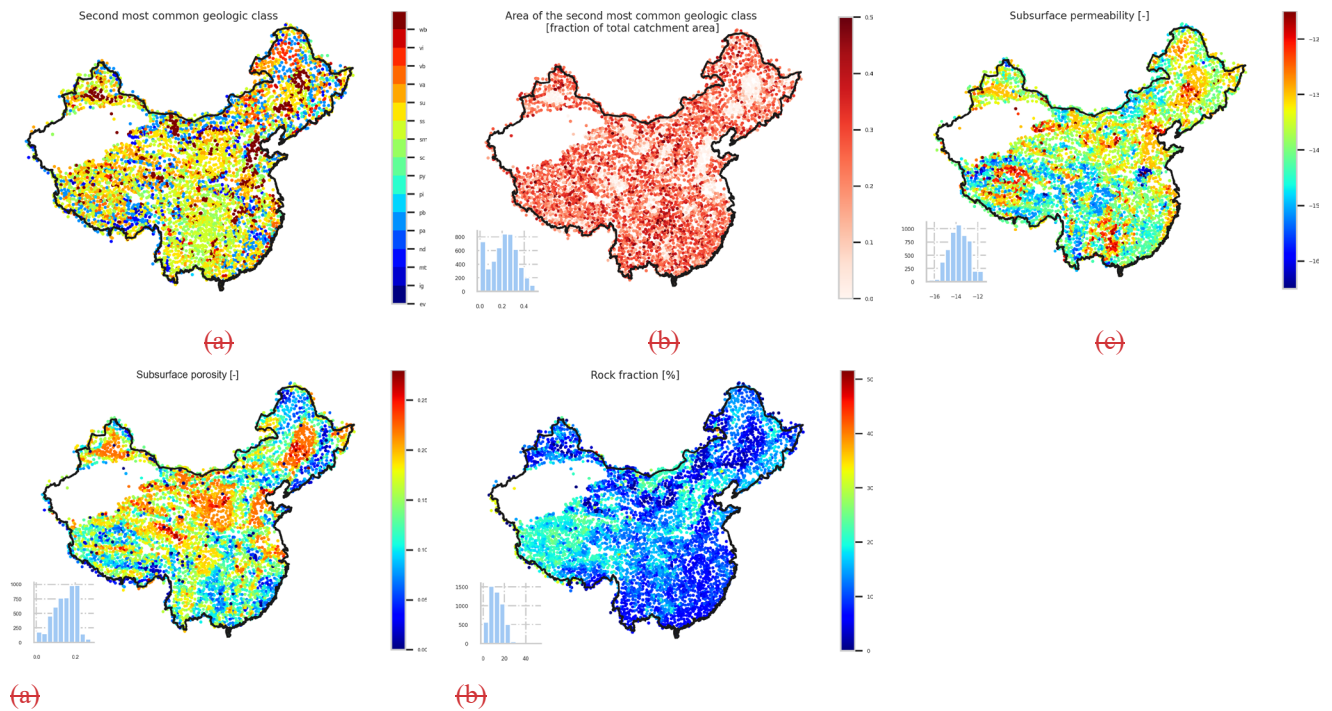
## 305 4 Geology

To describe the lithological characteristics of each catchment, we used the same two global datasets as CAMELS, Global Lithological Map (GLiMGLiM) (Hartmann and Moosdorf 2012) and GLobal HYdrogeologyGlobal Hydrogeology MaPS (GLHYMPS) (Gleeson, Moosdorf et al. 2014). Figure 4 presents the ~~results~~distributions of the geological types.

310 GLiM provides a high resolution global lithological map assembled from existing regional geological maps; it has been widely used for constructing datasets (e.g. SoilGrids250m (Hengl, Mendes de Jesus et al. 2017)). However, the data quality of GLiM can vary in different spatial locations depending on the quality of the original regional geological maps. GLiM consists of three levels, the first level contains 16 lithological classes, and the additional two levels describe more specific lithological characteristics. The GLiM is represented by 1,235,400 polygons; the polygons are converted to raster format for the basin-  
315 scale lithological type statistics. ~~For contiguous~~For China, the compiled regional data sources (China 1991, Xinjiang 1992, Survey 2001) have slightly lower resolutions than the GLiM target resolution (1:1 000 000). However, for a basin-scale study with a mean basin area of over 2000 km<sup>2</sup>, the classification accuracy should satisfy most applications.



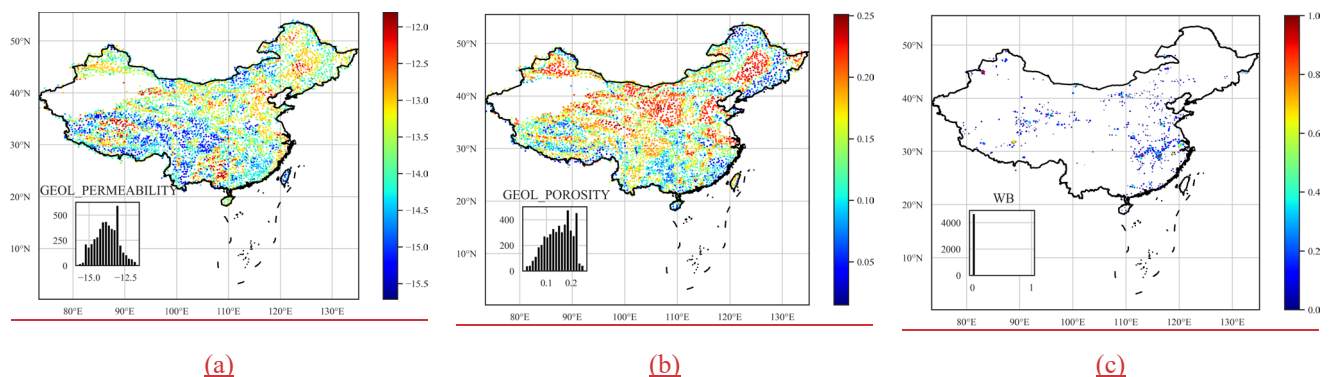


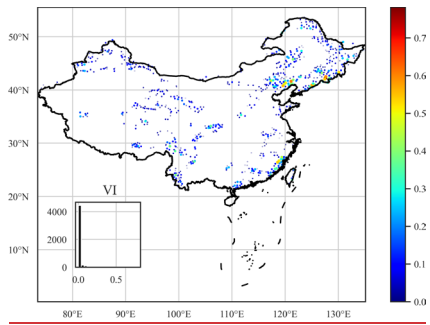


**Figure 4. Maps of geological characteristics over contiguous China. The histograms indicate the number of catchments (out of 4875) in each bin.**

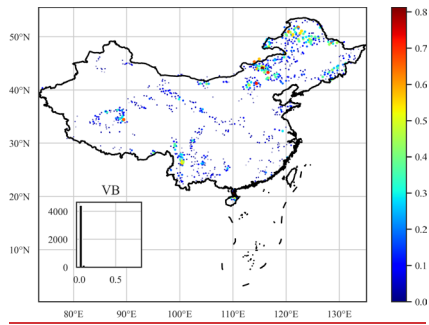
320 Compared to CAMELS and CAMELS-CL, one design consideration of the proposed dataset is that it should be more prepared for the data driven research, such that we aim to generate as many types of catchment scale data as possible since advanced data driven methods can learn the representation of inputs automatically. To this end, we determined and recorded Different from CAMELS and CAMELS-CL, we determined each lithological class's contribution to the catchment instead of recoding just the first and second most frequent classes. ~~The GLiM is represented by 1,235,400 polygons; the polygons are converted to raster format for the basin scale lithological type statistics.~~

325

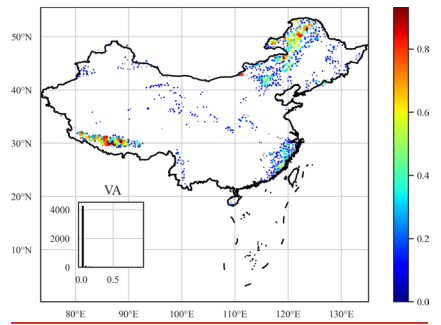




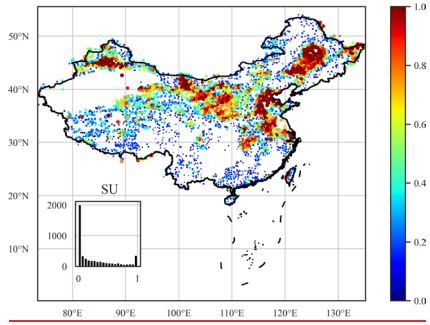
(a)



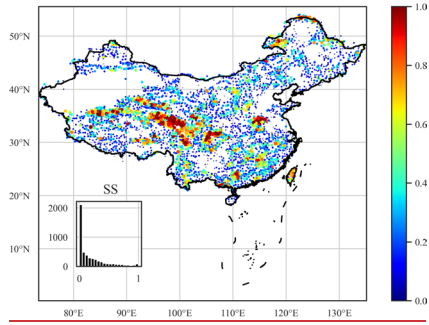
(b)



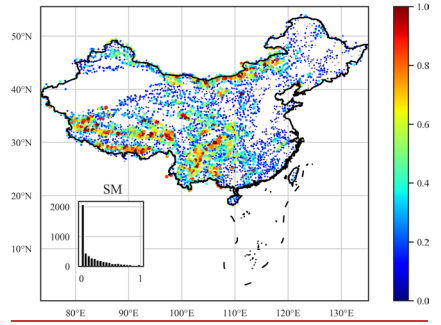
(c)



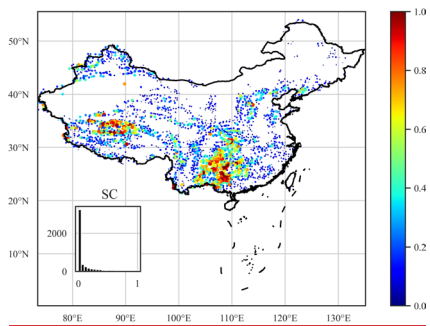
(d)



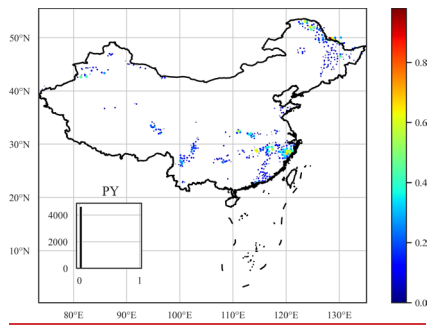
(e)



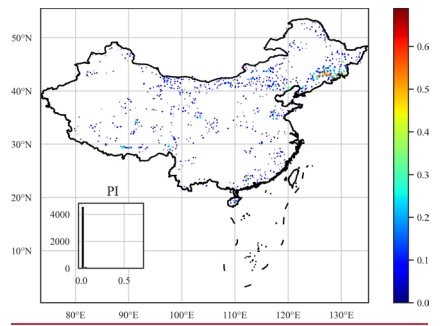
(f)



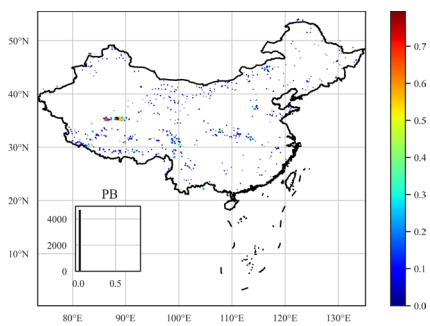
(g)



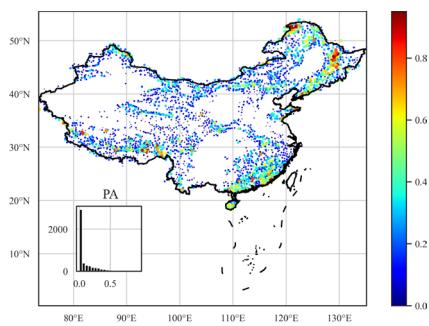
(h)



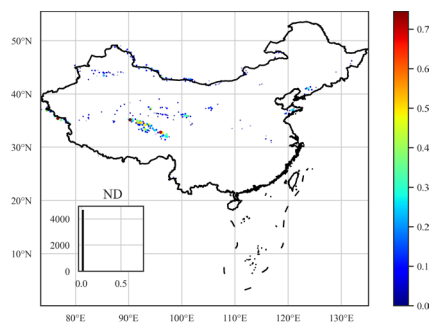
(i)



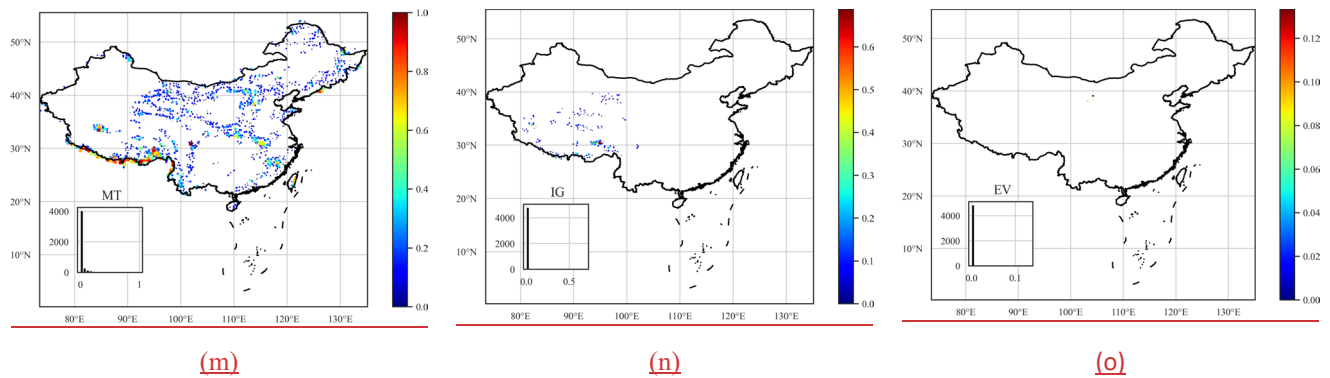
(j)



(k)



(l)



**Figure 4: Distributions of geological characteristics over China. For lithologies, the plot size is scaled by the lithology proportion.**

330 GLobal HYdrogeology MaPS (GLHYMPS) provides a global estimation of subsurface permeability and porosity, two critical characteristics for the soils' hydrological classification. Porosity and permeability influence an area's infiltration capacity. Soil with high porosity is likely to contain s amounts of water, and high permeable soil transmits water relatively quickly. Based on the high-resolution map of GLiM, which can differentiate fine and coarse-grained sediments and sedimentary rocks, GLHYMPS determined subsurface permeability depending on the different permeabilities of rock types. For the proposed dataset, we calculated the catchment arithmetic mean for porosity. Followed (Gleeson, Smith et al. 2011), the logarithmic scale geometric mean is used for representing subsurface permeability. The summary of geological characteristics is present in Table 3.

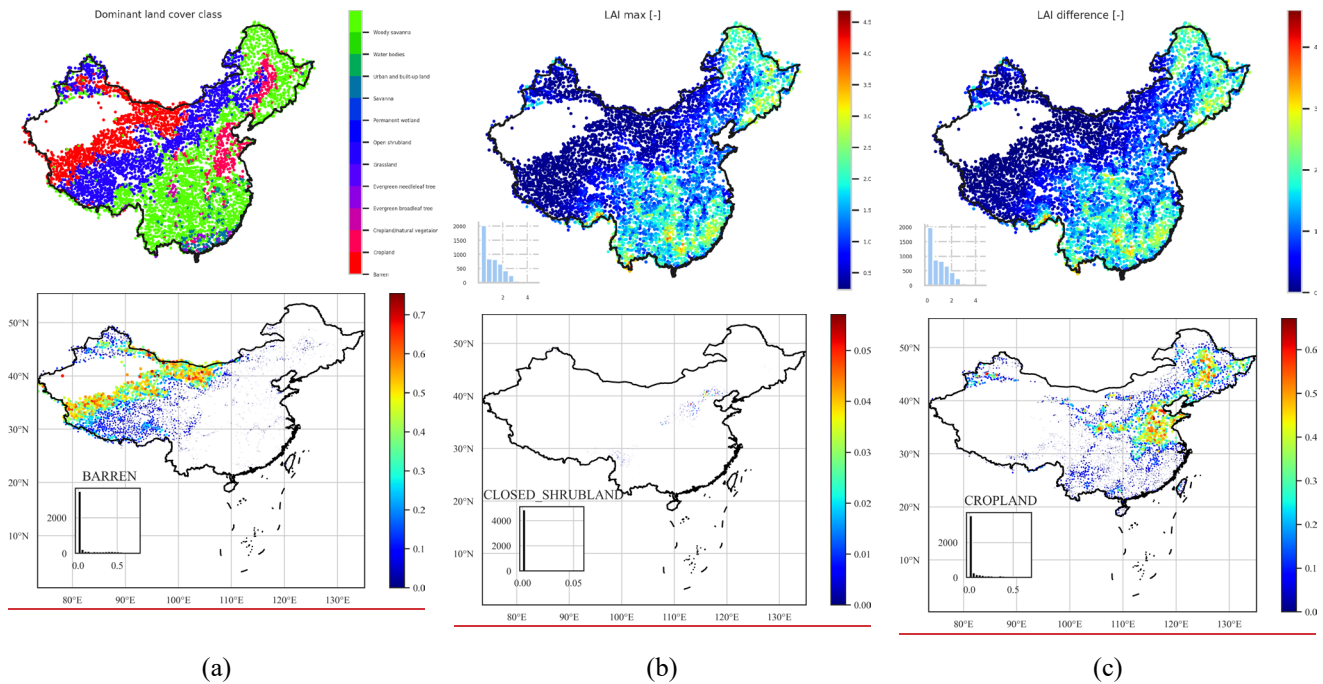
340 Porosity and permeability have similar distributions as geological classes. These two characteristics are highly dependent on rock properties, unconsolidated sediments, mixed sedimentary rocks, siliciclastic sedimentary rocks, carbonate sedimentary rocks, and acid plutonic rocks are the five most common geological classes in ~~contiguous~~-China. Unconsolidated sediment is the most common rock type in ~~contiguous~~ China, dominating 31.9% of catchments; it extends from Xinjiang to the inland of the northeast and the coastal area surrounding the Bohai Sea, due to the high proportion of unconsolidated sediments present in the rock, these areas typically have high permeability and medium porosity. Mixed sedimentary rocks are the second most common rock type in ~~contiguous~~ China, accounting for 20.3% of catchments, it dominated the southern Qinghai-Tibet Plateau, western Yunnan-Guizhou Plateau, and northern Inner Mongolia. These areas typically have high porosity and low permeability. Siliciclastic sedimentary rocks dominate 17.7% of basins, mainly distributed in the northern part of the Qinghai-Tibet Plateau and the junction of the Qinghai-Tibet Plateau and the Yunnan-Guizhou Plateau; there are also some distributions in the eastern inland. These areas have low subsurface permeability and high subsurface porosity. Amongst all catchments, 9.8% of catchments are dominated by carbonate sedimentary rocks. Carbonate sedimentary rocks are mainly located in eastern Yunnan and northern Qinghai-Tibet Plateau. Acid plutonic rocks are typically distributed in the mountains surrounding the inland northeast, namely the Daxinganling Mountain and the hills in southern Guangdong and southwestern Guangxi. They are also distributed along the Brahmaputra river in the south part of the Qinghai-Tibet Plateau. The distribution of Acid plutonic

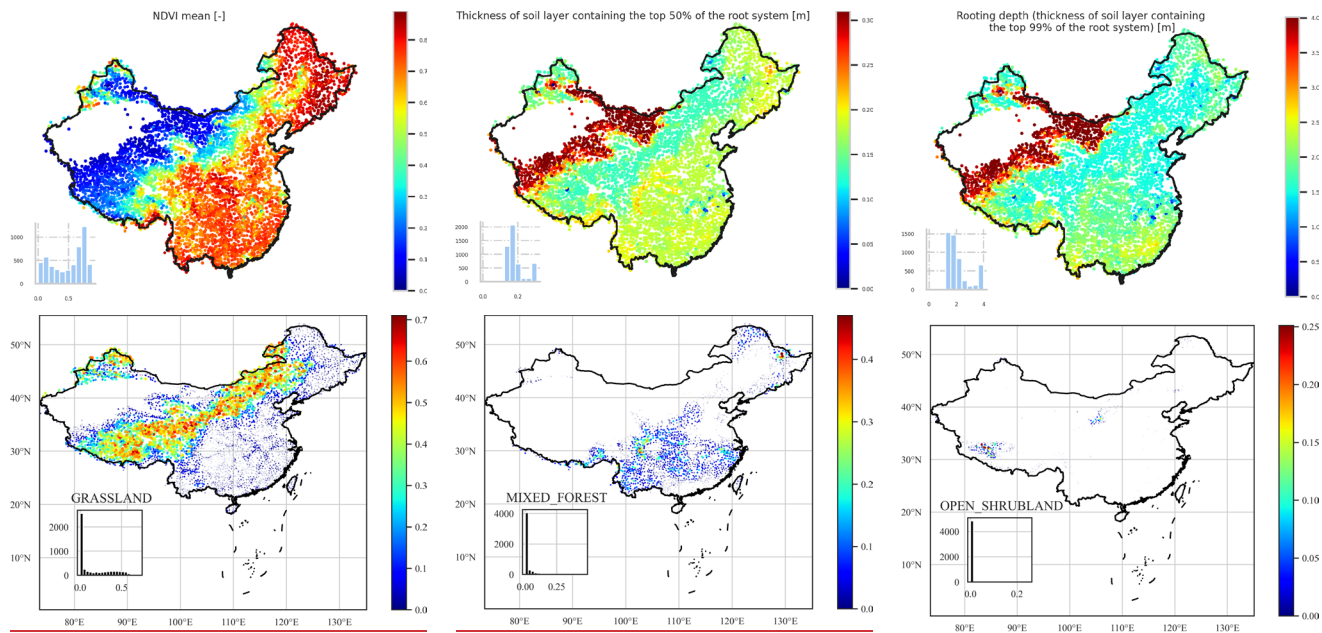


rocks is relatively scattered; there are many isolated Acid plutonic rocks distributions in different locations of **contiguous** China, accompanied by medium permeability and high porosity.

355 **In summary, the** types of rocks in **contiguous** China are dominated by unconsolidated sediments and mixed sedimentary rocks. In 33.86% of the catchments, the dominant rock types occupy less than 50% of the catchment areas, and only 16.8% **of** basins are having a dominant rock type with an area fraction greater than 90%. Amongst **48754911** basins, 9.4% of basins have prevalent rock types wholly occupying the area.

## 5 Landcover

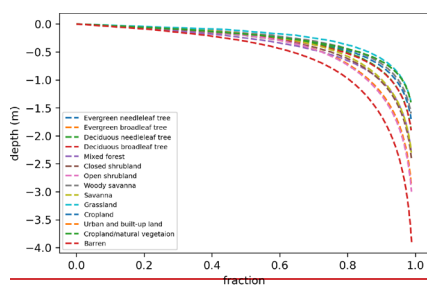




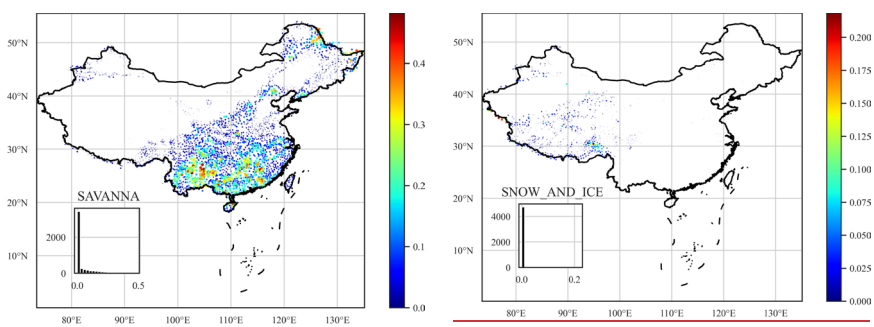
(d)

(e)

(f)

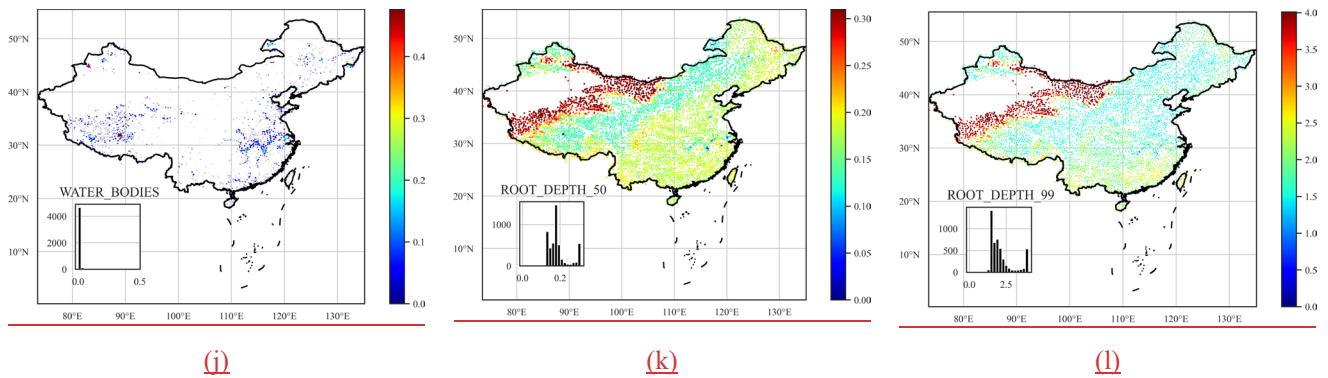


(g)



(h)

(i)



(j)

(k)

(l)

360 **Figure 5. Maps: Distributions of land cover characteristics over contiguous China. The histograms indicate For land cover types, the number plot size is scaled by the size of catchments (out of 4875) in each bin. the land cover proportion.**

We selected two indicators to characterize vegetation density and growth on the surface: Normalized difference vegetation index (NDVI) and Leaf area index (LAI). NDVI is an indicator with a valid range of -0.2 to 1, assessing whether the area being observed contains live green vegetation or the plants' health. However, NDVI is just a qualitative measurement of the vegetation density; it cannot provide a quantitative estimate of the vegetation density in the area. Moreover, NDVI often provides inaccurate vegetation density measurements, and only long-term measurement and comparison can ensure its accuracy. NDVI alone is not enough to estimate the state of plants in an area. Therefore, we have selected another indicator, LAI, to supplement the deficiencies of NDVI.

370 LAI is defined as the total needle surface area per unit ground area and half of the entire needle surface area per unit ground surface area. It is a quantifiable value. It is functionally related to many hydrological processes like water interception (van Wijk and Williams 2005). (Buermann, Dong et al. 2001) verifies the validity of LAI used to characterize vegetation growth. The data sources used are The Terra Moderate Resolution Imaging Spectroradiometer (MODIS) Vegetation Indices (Didan 2015) for NDVI and Moderate Resolution Imaging Spectroradiometer (MODIS) (Myneni, Knyazikhin et al. 2015) for LAI. 375 Followed (Addor, Newman et al. 2017), we determined maximum monthly LAI as an indicator characterising vegetation interception capacity and the maximum evaporative capacity and the difference between the maximum and minimum monthly LAI representing LAI's temporal variations.

Land cover classification refers to segmenting the ground into different categories based on remote sensing images. The Terra and Aqua combined Moderate Resolution Imaging Spectroradiometer (MODIS) Land Cover Type provides different results depending on the classification system used. Annual International Geosphere-Biosphere Programme (IGBP) classification is used for building the dataset, which is derived by the c4.5 decision tree algorithm. The IGBP classification system was formulated by the IGBP Land Cover Working Group in 1995, resulting in 17 categories of land cover types (Belward, Estes et al. 1999). (Friedl, Sulla-Menashe et al. 2010) compared the IGBP data of MODIS with other reference dataset Friedl, Sulla-Menashe et al. (2010) compared the IGBP data of MODIS with other reference datasets and concluded that the MODIS

classification of IGBP has an accuracy of 75%. We determined the fraction of each land cover class for each basin based on the Terra and Aqua combined Moderate Resolution Imaging Spectroradiometer (MODIS) Land Cover Type (Sulla-Menashe and Friedl 2018), which differentiates our dataset from CAMELS and CAMELS-CL (only calculated the proportion of the dominant types).

390

Followed (Addor, Newman et al. 2017), we ~~also~~ computed the average rooting depth (50% and 90%) for each catchment based on the IGBP classification using a two-parameter method (Zeng 2001). The root depth distribution of vegetation affects the ground's water holding capacity and the topsoil layer's annual evapotranspiration (Desborough 1997). Many models use root depth as an essential parameter to characterize soil moisture absorption capacity. (Zeng 2001) developed a two-parameter asymptotic equation for estimating root depth distribution; the root depth distribution is global, derived based on the IGBP classification avoiding the problem of significantly different root distributions in various research. Figure 5(g) shows root depth distributions of different vegetation types, based on (Zeng 2001)'s ~~method~~. The 90% root depth is usually considered to be "rooting depth", among the 17 categories of IGBP, cropland has the smallest rooting depth, and open shrubland has the largest. The 90% root depth of all vegetation is less than 2 meters. The national distribution of catchments soil characteristics is shown in Fig. 5.

395

400

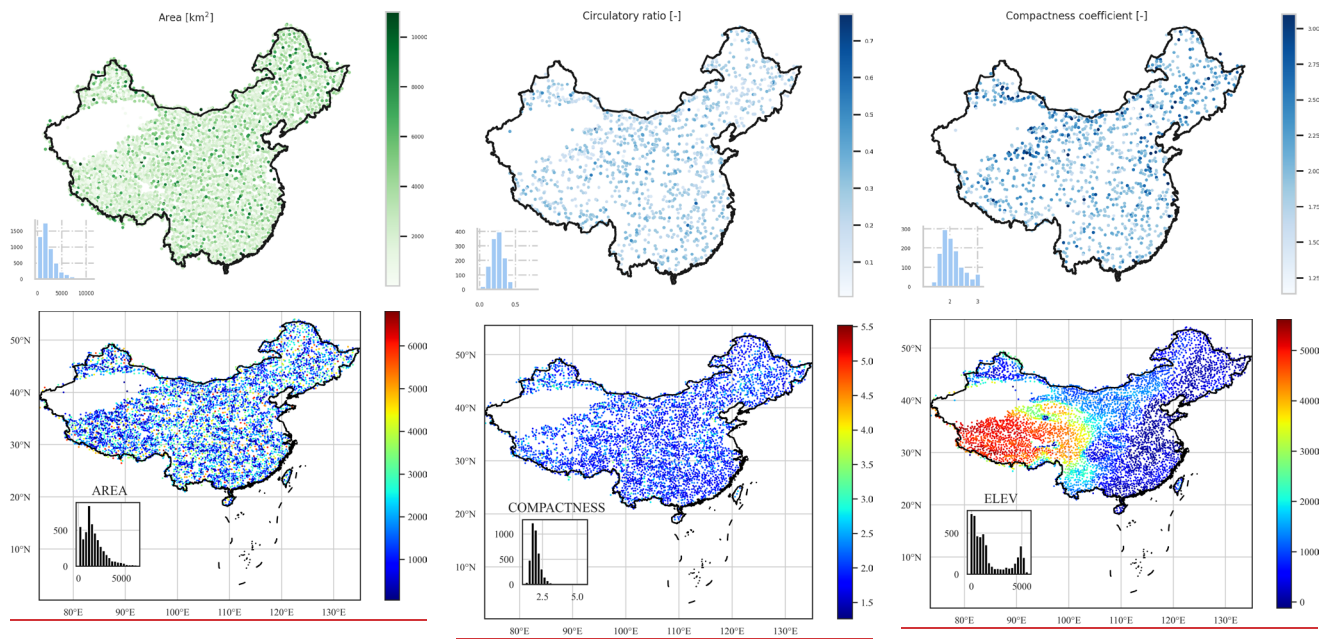
## 6 Location and topography

The catchments' boundary files are obtained from the global drainage basin dataset (Masutomi, Inui et al. 2009). The ~~PDBD~~GDBD dataset was derived from digital elevation models (DEMs) with a high-resolution (100m-1km), and the errors were corrected by either automatic methods or manually. Additionally, ~~PDBD~~GDBD also provides population and population density estimates for catchments, and these two indicators are also included in our dataset as a measure of human intervention. Global ~~Runoff~~Streamflow Data Centre (Center 2005) discharge gauging stations were used for referencing the derived basins. ~~In contiguous China, PDBD~~GDBD has a high average match area rate (AMAR) and good geographic agreement with existing global drainage basin data in China. Based on the high-quality dataset, precise geographic and topographic information can be derived. ~~See Fig. 6 for a summary.~~

405

410

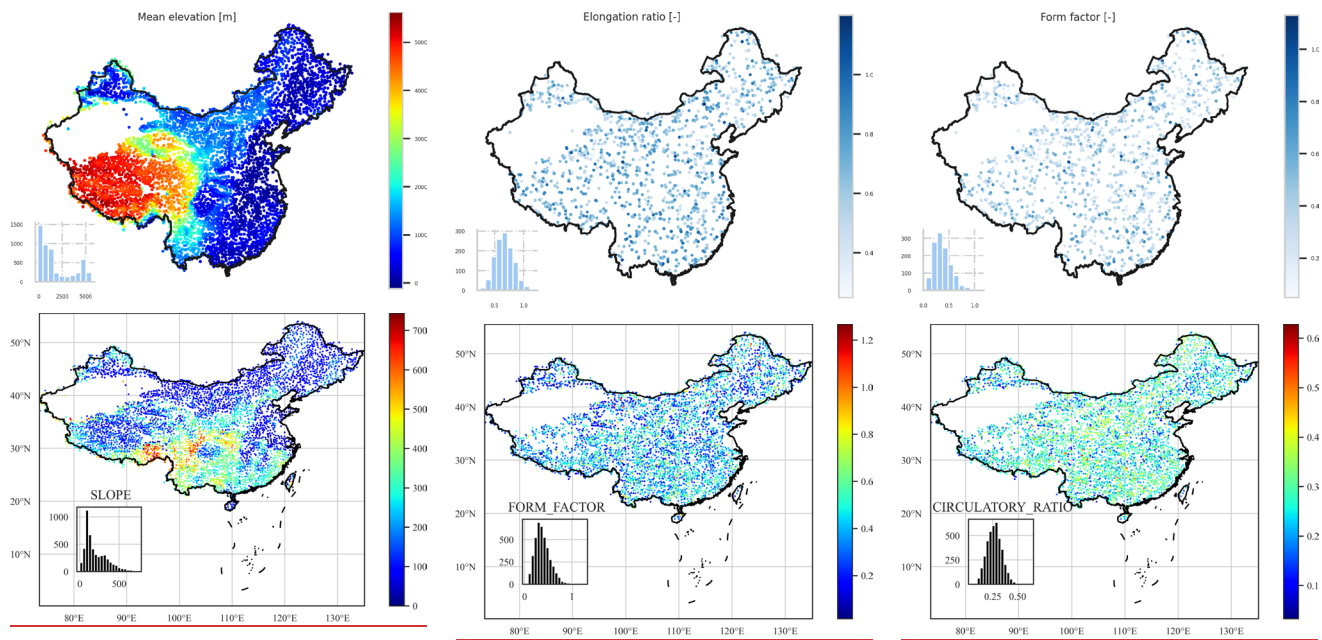
The topography attributes of each catchment are determined based on the ASTGTM product retrieved from <https://lpdaac.usgs.gov>, maintained by the NASA EOSDIS Land Processes Distributed Active Archive Center (LP DAAC) at the USGS Earth Resources Observation and Science (EROS) Center.



(a)

(b)

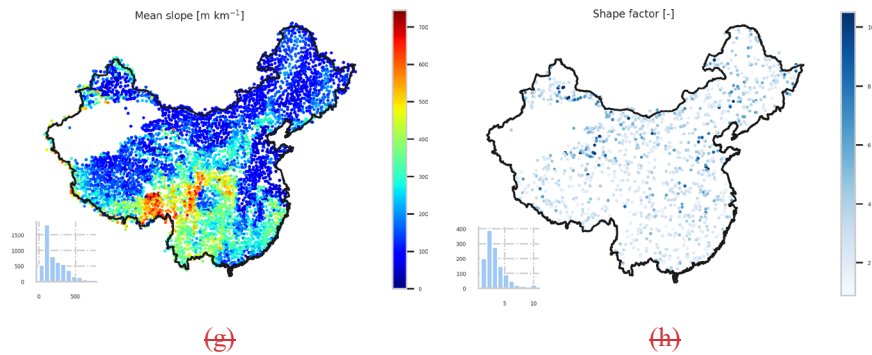
(c)



(d)

(e)

(f)



415 **Figure 6. ~~Maps~~Distributions of topographic characteristics over contiguous China. The histograms indicate the number of catchments (out of 4875) in each bin.**

The CAMELS dataset ~~just~~ provides two parameters (two area estimates) for describing the catchment shape; ~~however, the~~ The physical characteristics of a catchment can affect the ~~runoffstreamflow~~ volume and the ~~runoffstreamflow~~ hydrograph of the catchment under a storm. To provide a complete description of the catchment shape, we computed several geometrical parameters of the catchment related to the ~~runoffstreamflow~~ process, (Fig. 6), including catchment form factor, shape factor, compactness coefficient, circulatory ratio and the elongation ratio (Subramanya 2013). A summary of the location and topography attributes can be found in Table 3.

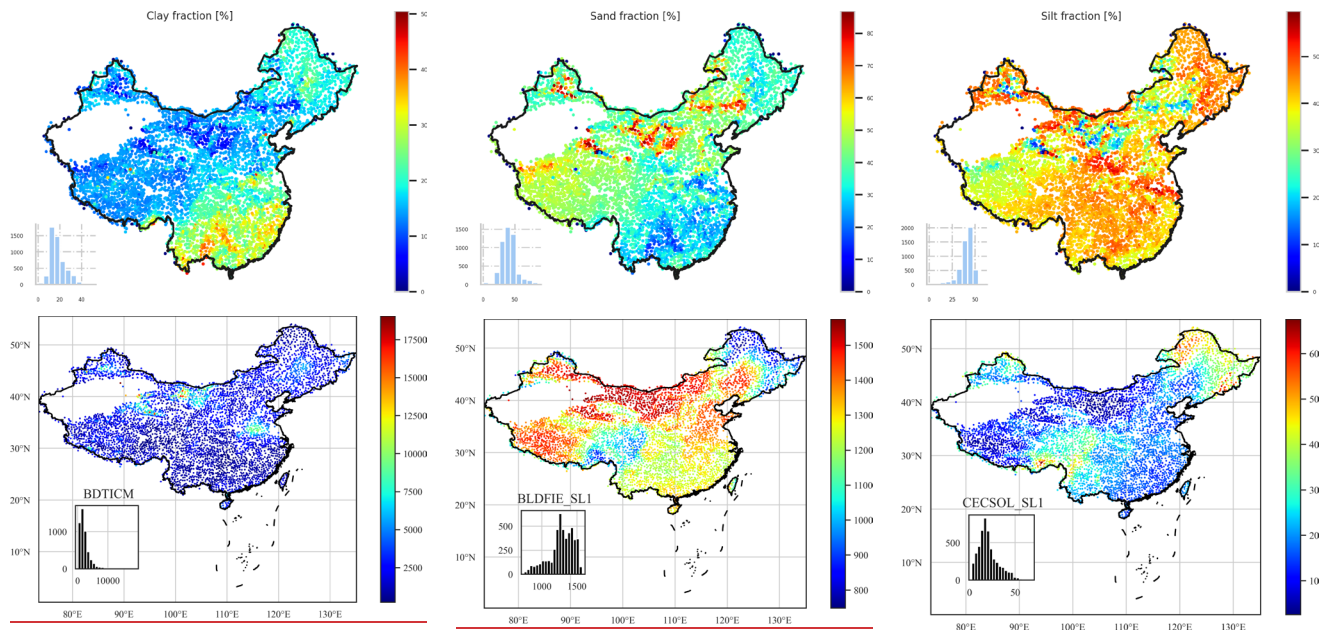
420

## 7 Soil

The proposed dataset has a total of 54 soil attributes (Table 3) derived from (Hengl, Mendes de Jesus et al. 2017), (Dai, Xin et al. 2019) and (Shangguan, Dai et al. 2013). ~~The summary result is shown in Fig. 7.~~ Five categories of soil characteristics (pH in H<sub>2</sub>O, organic carbon content, depth to bedrock, cation-exchange capacity, and bulk density) are determined from SoilGrids. SoilGrids (Hengl, Mendes de Jesus et al. 2017) provides global predictions for soil properties including organic carbon, bulk density, cation exchange capacity (CEC), pH, soil texture fractions and coarse fragments by fusing multiple data sources including MODIS land products, SRTM DEM, climatic images and global landform and lithology maps at the 250m resolution. (Fig. 7). SoilGrids made predictions based on machine learning algorithms and many ~~eo~~covariates layers primarily derived from remote sensing data. SoilGrids has soil characteristics for several soil depths.

430

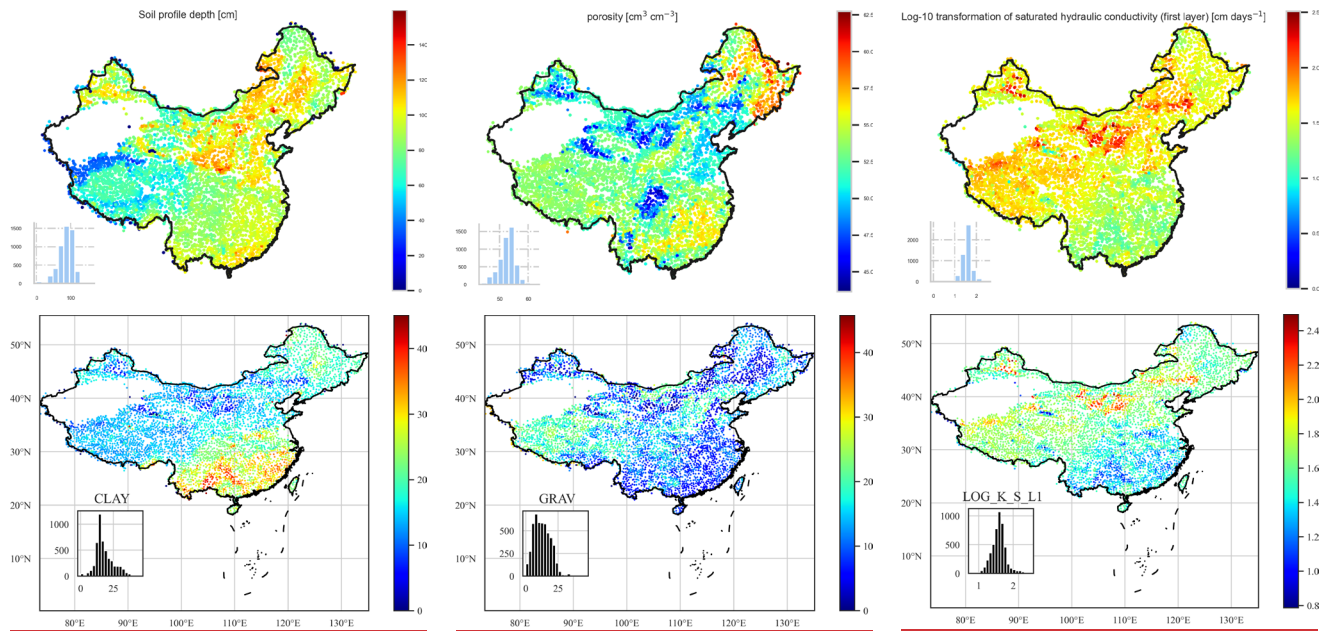




(a)

(b)

(c)

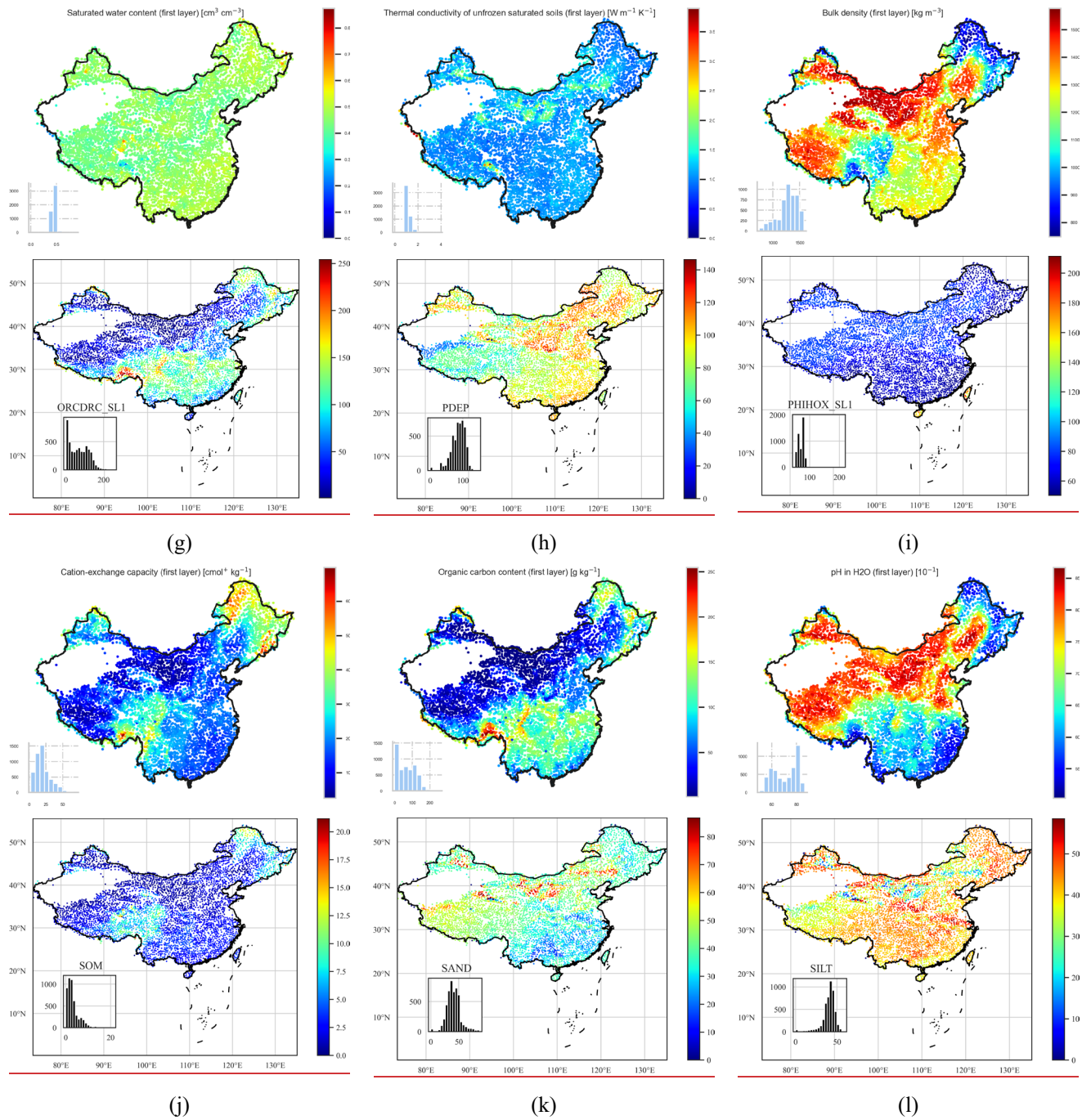


(d)

(e)

(f)





**Figure 7. Maps: Distributions of soil characteristics over contiguous China. The histograms indicate the number of catchments (out of 4875) in each bin.**

435 ~~Unlike~~Different from CAMELS, whose reported results are obtained by a linear weighted combination of the different soil layers, and CAMELS-BR, whose products are soil characteristics at a depth of 30cm. We computed soil characteristics at all soil layers provided by ~~SoilGrids such that advanced models can learn directly from the raw inputs~~SoilGrids250m.

~~To be consistent with CAMELS, we also~~We determined saturated water content and saturated hydraulic conductivity (Dai, 440 Xin et al. 2019). ~~We also introduced thermal conductivity of unfrozen saturated soils (Dai, Xin et al. 2019). (Dai, Xin et al. 2019)~~Based on the same dataset, we also introduced the thermal conductivity of unfrozen saturated soils. Dai, Xin et al. (2019) provides a global estimation of soil hydraulic and thermal parameters using multiple Pedotransfer Functions (PTFs) based on ~~SoilGrids~~the SoilGrids250m dataset. Based on the ~~SoilGrids~~SoilGrids250m and GSDE (Shangguan, Dai et al. 2014) datasets, ~~(Dai, Xin et al. 2019) produced six soil layers with a spatial resolution~~Dai, Xin et al. (2019) produced six soil layers with a 445 spatial resolution of 30×30 arc-second. The vertical resolution of (Dai, Xin et al. 2019) is the same as the ~~SoilGrids~~SoilGrids250m, with six intervals of 0–0.05 m, 0.05–0.15 m, 0.15–0.30 m, 0.30–0.60 m, 0.60–1.00 m, and 1.00–2.00 m. ~~Same as the methods applied to SoilGrids, we determined~~We determine and ~~records~~record catchment soil characteristics for all these layers.

450 ~~To provide even more complete description of the soil~~In addition, we determined seven more soil characteristics (Shangguan, Dai et al. 2013) including soil profile depth, porosity, clay/silt/sand content, rock fragment, and soil organic carbon content. ~~(Shangguan, Dai et al. 2013)~~Shangguan, Dai et al. (2013) provides physical and chemical attributes of soils derived from 8979 soil profiles at 30×30 arc-second resolution, the polygon linkage method was used to derive the spatial distribution of soil properties. The profile attribute database and soil map are linked under a framework avoiding uncertainty in taxon referencing.

455 Depth to bedrock controls many physical and chemical processes in soil. The distribution of depth to bedrock in ~~contiguous~~ China is characterised by (i) low in the mountainous areas, such as Yunnan province and Chongqing City; (ii) high in barren areas, e.g. North and Northwest China. The introduced soil pH value is crucial since it influences many other physical and chemical soil characteristics. The spatial variability of soil pH in ~~contiguous~~ China is characterised by (i) soils in southern 460 ~~contiguous~~ China are acid to strongly acid; (ii) soils in northern China are natural or alkaline; (iii) soils in ~~north-eastern~~northeastern forested areas are also acid (pH < 7.2). Cation exchange capacity can be seen as a measure of soil fertility since it measures how much nutrient the soil can store such that it influences the growth of the ~~vegetations~~vegetation. Cation exchange capacity is positively correlated with soil organic matter content and clay content, which Cation exchange capacity is generally low in sandy and silty soils. The spatial variability of Cation exchange capacity in ~~contiguous~~ China is 465 characterised by (i) high in peat and forested areas in Qinghai-Tibet Plateau, central and northeast China (ii) The Cation exchange capacity in the desert area such as the northwest is extremely low. Soil hydraulic and thermal properties are greatly affected by soil organic matter (SOM). Soil organic matter has a similar distribution to the cation exchange capacity: high in the peat and forested areas such as northeast China and low in the north and northwest.

## 8 Meteorological time series

470 **Table 4:** Summary table of catchment meteorological time series available in the proposed dataset

Variable	Description	Unit
prs	catchment daily averaged ground pressure	hPa
tem	catchment daily averaged temperature at 2 m above ground	°C
rhu	catchment daily averaged relative humidity	-
pre	catchment daily averaged precipitation	mm d <sup>-1</sup>
evp	catchment daily averaged evaporation measured by ground instruments	mm d <sup>-1</sup>
win	catchment daily averaged wind speed at 2 m above ground	m s <sup>-1</sup>
ssd	catchment daily averaged sunshine duration	h d <sup>-1</sup>
gst	catchment daily averaged ground surface temperature	°C
pet	catchment daily averaged potential evapotranspiration determined by Penman's equation (see Appendix A)	mm d <sup>-1</sup>

There have been many studies based on SURF\_CLI\_CHN\_MUL\_DAY in China (Liu, Xu et al. 2004, Xu, Gao et al. 2009, Huang, Han et al. 2016, Liu, Zheng et al. 2017), such as trend analysis of the pan evaporation (Liu, Yang et al. 2010). Still, there has not yet been a large-scale basin-oriented meteorological time series dataset in ~~contiguous~~ China. Researchers still need to do repeated works to extract historical meteorological data from the SURF\_CLI\_CHN\_MUL\_DAY dataset for the research. For the first time, we release a catchment-scale meteorological time series dataset. ~~We will also~~The open-source ~~thesourced~~ code ~~for researchers to can~~ generate any catchment's meteorological time series within ~~contiguous~~ China. The basin-oriented dataset provides meteorological time series for ~~48754911~~ basins from 1990 to ~~20182020~~ based on the China Meteorological Data Network. Meteorological time series includes pressure, temperature, relative humidity, precipitation, evaporation, wind speed, sunshine duration, ground surface temperature and potential evapotranspiration (~~see (Table 4 for a summary).~~).

The meteorological time series data from 1951 to 2010 is derived based on the "1951-2010 China National Ground Station Data Corrected Monthly Data File Basic Data Collection" data construction project. Other data include monthly reported data to the National Meteorological Information Centre by the provinces, and hourly and daily data uploaded by automatic ground stations in real-time. ~~The SURF\_CLI\_CHN\_MUL\_DAY dataset is quality controlled, the quality and completeness of each variable are significantly improved compared to the previous similar products. M~~During the development of the dataset, missing data were filled by interpolating its nearest stations.

490 Figure 2 presents the variation of the ~~distribution number~~ of ~~the observation~~-sites. The start date of the recording is 1951, but  
because the early site distribution is sparse, we only used records from 1990 to ~~2018~~2020 to construct the dataset to ensure the  
data quality. ~~The interpolation method used is the~~Inverse distance weighting ~~since it~~ shows better performance than other  
~~comparators. Catchment scale raster is extracted from the interpolated national raster using the open-source rasterio.<sup>6</sup> package.~~  
495 ~~For all variables, we take the arithmetic mean on the extracted catchment raster as the catchment mean. Potential~~interpolation  
methods. In addition, potential evapotranspiration (PET) is estimated based on Penman's Equation (Appendix A) and other  
~~catchment~~ meteorological variables.

### ~~9 Normal Camels YR—Normalized Catchment attributes and meteorology for Yellow River basin~~

~~Apart from the dataset providing the catchment attributes and meteorological forcing for contiguous China, we also offer a~~  
~~self-contained dataset covering the Yellow River basin with normalized streamflow measurements. The streamflow data are~~  
500 ~~normalized to have zero mean and a standard deviation of 1 for each basin. The Normal Camels YR dataset is designed to~~  
~~support machine learning and deep learning research related to hydrology. In particular, fifty-four watersheds are less affected~~  
~~by human activities (selection is based on the Global Reservoirs and Dam databases (GRanD) (Lehner, Liermann et al. 2011)~~  
~~which provides the locations of reservoirs and dams globally), which makes them suitable for rainfall runoff modelling~~  
~~research. For most machine learning and deep learning algorithms, data normalization will not affect model performance (e.g.,~~  
505 ~~neural network-based and tree-based algorithms). Besides, other research, such as trend analysis, can also be carried out. The~~  
~~Normal Camels YR dataset is self-contained to fully describe the Yellow River basin and is particularly helpful for the~~  
~~hydrology research of the Yellow River.~~

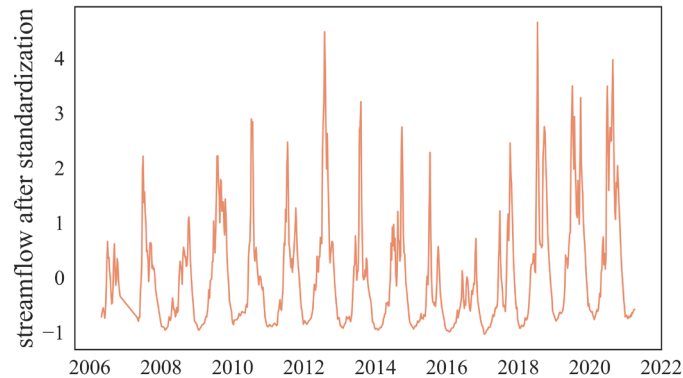
~~During the dataset development, basins with too few observations are removed, resulting in discontinuous basin identifiers.~~  
510 ~~Normal Camels YR covers 102 gauges in the Yellow River basin, providing basin boundary shapefiles, static attributes and~~  
~~normalized streamflow measurements for each basin. The covered basins have areas ranging from 134 to 804,421 square~~  
~~kilometres. The time resolution of streamflow measurements is seven days, and the mean length of records of the streamflow~~  
~~measurements is 684, which means the mean period of the streamflow measurements for each basin is over 13 years.~~  
~~Meteorological variables included in Normal Camels YR is slightly different; it introduced daily maximum and minimum for~~  
515 ~~some variables (Table 5).~~

---

<sup>6</sup> <https://github.com/mapbox/rasterio>

## 9 HydroMLYR: Hydrology dataset for Machine Learning in YRB

In addition to the basin-wise static attributes provided in CCAM, we propose HydroMLYR, a hydrology dataset for machine learning research in the YRB (Fig. 1). HydroMLYR includes standardized streamflow measurements for 102 basins. The streamflow data is seven-day averaged and standardized basin-wise to have zero mean and a standard deviation of 1 (Fig. 8). The HydroMLYR dataset is proposed to support machine learning or deep learning hydrology research (e.g., neural network-based and tree-based algorithms). It can be used in two cases: (1) to develop machine learning models on the YRB or (2) when it is desirable to verify the generalization ability of a machine learning model on YRB.



**Figure 8: Examples of standardized runoff**

The dataset provides 40 natural basins in the dataset which are not affected by reservoirs and dams. The selection is based on a newer version<sup>7</sup> of the Global Reservoirs and Dam databases (Lehner, Liermann et al. 2011) which provides the locations of reservoirs and dams globally. HydroMLYR covers 102 basins in the YRB, including basin boundary shapefiles, static attributes, and standardized streamflow measurements for each basin. The covered basins have areas ranging from 134 to 804,421 square kilometres. Therefore, modelling on a large scale of the YRB is also possible. Meteorological records in HydroMLYR introduced daily maximum and minimum for some forcing variables (Table 5).

The original streamflow observations are not continuous. The average record length is 11.3 years. Although the development of machine learning models does not necessarily require the data to be continuous, we separately provide continuous streamflow observations with an average record length of 8.3 years.

**Table 5: Meteorological variables provided in Normal-Camels-YR, the time series length is 22 years (1999-2020) HydroMLYR**

Attribute name	Description	Unit
evp	catchment daily averaged evaporation (observations)	0.1 mm d <sup>-1</sup>
gst_mean	catchment daily averaged ground surface temperature	0.1 °C

<sup>7</sup> [http://globaldamwatch.org/data/#core\\_global](http://globaldamwatch.org/data/#core_global)

gst_min	catchment daily minimum ground surface temperature	0.1°C
gst_max	catchment daily maximum ground surface temperature	0.1°C
pre	catchment daily averaged precipitation	0.1 mm d <sup>-1</sup>
prs_mean	catchment daily averaged ground surface pressure	0.1 hPa
prs_max	catchment daily maximum ground surface pressure	0.1 hPa
prs_min	catchment daily minimum ground surface pressure	0.1 hPa
rhu	catchment daily averaged relative humidity	-
ssd	catchment daily averaged sunshine duration	0.1 h
tem_mean	catchment daily averaged temperature	0.1°C
tem_min	catchment daily minimum temperature	0.1°C
tem_max	catchment daily maximum temperature	0.1°C
win_max	catchment daily maximum wind speed	0.1 m s <sup>-1</sup>
win_mean	catchment daily averaged wind speed	0.1 m s <sup>-1</sup>

## 10 Data ~~and code~~ availability ~~and software packages used.~~

540 The proposed dataset is freely available at <http://doi.org/10.5281/zenodo.47040175137288>. The files provided are (i) several separate files containing 120+ catchments attributes, (ii) the daily meteorological time series in a zip file, (iii) the catchment boundaries used to compute the attributes and extract the time series, (iv) the ~~Normal Camels YRHydroMLYR~~ dataset, (v) an attribute description file and (v) a readme file. ~~The code used to generate the dataset is mainly based on several publicly available packages: rasterio, gdal<sup>8</sup>, pyshp<sup>9</sup>, geopandas<sup>10</sup>, fiona<sup>11</sup>, and xarray<sup>12</sup>. Complement code for generating any watershed's dataset will be released soon.~~

## 11 Conclusion

545 The CCAM dataset proposed in this paper provides a novel dataset for hydrological research in ~~contiguous China. In the study China area, there is no catchment attributes dataset has been proposed before, either a catchment scale time series~~

<sup>8</sup> <https://github.com/OSGeo/gdal>

<sup>9</sup> <https://github.com/GeospatialPython/pyshp>

<sup>10</sup> <https://github.com/geopandas/geopandas>

<sup>11</sup> <https://github.com/Toblerity/Fiona>

<sup>12</sup> <https://github.com/pydata/xarray>

~~meteorological dataset.~~ All ~~catchments~~basins delaminated from the DEM are studied, covering ~~contiguous~~entire China. The dataset includes daily meteorological forcing time-series data including precipitation, temperature, potential evapotranspiration, wind, ground surface temperature, pressure, humidity, sunshine duration and derived potential evapotranspiration of ~~48754911~~ catchments. The proposed time series dataset is derived based on the quality-controlled ~~site observation dataset,~~ SURF\_CLI\_CHN\_MUL\_DAY. ~~We will also release the complement code for generating any shapefile's meteorological time series within contiguous China based on the SURF\_CLI\_CHN\_MUL\_DAY dataset (freely available for Chinese researchers).~~ The dataset has longer time series (from 1990 to 2018) and more meteorological variables than the ~~previously proposed datasets.~~ The dataset also ~~dataset.~~ CCAM includes 120+ catchment attributes, including soil, land cover, geology, climate indices and topography for each catchment. ~~We produced a series of maps depicting the catchment attributes distributions in contiguous China. These maps present regional changes of various features; we also describe~~estimate the relationships between them. ~~The integration of~~ based on Kendall's correlation. Integrating multiple data sources into one dataset at a catchment-scale ~~dramatically~~ simplifies the data compilation process in research. ~~Based on the dataset, we~~ CCAM can help test hypotheses and formulate valid conclusions under various conditions, not just limited to a few specific locations. ~~Together with the Normal Camels YR dataset, the proposed dataset can~~ and help explore how different basin characteristics influence hydrological behaviours, learn the migration of hydrological behaviours between different basins, and ~~to~~ develop general frameworks for large-scale model evaluation and benchmarking in China. A limitation of the study is the lack of estimation of the uncertainty of the meteorological time series. An alternative is to evaluate the uncertainty of the basin-wise meteorological data based on multiple independent data sources, but there are few data that provide as many data types as SURF\_CLI\_CHN\_MUL\_DAY. Hence, it poses a challenge for evaluating the uncertainty of these eight meteorological variables, which is left for future studies.

## Appendix A: Modified Penman's equation

Penman's equation (Subramanya 2013), incorporating some modifications to the original formula, is:

$$PET = \frac{AH_n + E_a\gamma}{A + \gamma}$$

where  $PET$  is the daily potential evapotranspiration in mm per day;  $A$  is the slope of the saturation vapour pressure ( $e_w$ ) vs temperature ( $t$ ) curve at the mean air temperature, in mm of mercury per Celsius;  $H_n$  is the net radiation in mm of evaporable water per day;  $E_a$  is a parameter including wind speed and saturation deficit;  $\gamma$  is the psychrometric constant = 0.49 mm of mercury per Celsius.

The relationship between  $e_w$  and  $t$  is defined as:

$$e_w = 4.584 \exp\left(\frac{17.27t}{237.3 + t}\right)$$

The following equation estimates the net radiation:



$$H_n = H_a(1 - r) \left( a + b \frac{n}{N} \right) - \sigma T_a^4 (0.56 - 0.092 \sqrt{e_a}) \left( 0.10 + 0.90 \frac{n}{N} \right)$$

580 where  $H_a$  is the incident solar radiation outside the atmosphere on a horizontal surface, expressed in mm of evaporable water  
per day (a function of the latitude and period of the year as indicated in Table A1);  $a$  is a constant depending upon the latitude  
585  $\phi$  and is given by  $a = 0.29 \cos \phi$ ;  $b$  is a constant = 0.52;  $n$  is the sunshine duration in hours;  $N$  is the maximum possible  
hours of bright sunshine (a function of latitude, see Table A2);  $r$  is the reflection coefficient;  $\sigma$  is the Stefan-Boltzman constant  
=  $2.01 \times 10^{-9}$  mm/day;  $T_a$  is the mean air temperature in degrees kelvin;  $e_a$  is the actual mean vapour pressure in the air in  
mm of mercury.

**Table A1:- Mean Monthly Solar Radiation,  $H_a$  in mm of Evaporable Water/Day**

North latitude	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
0°	14.5	15.0	15.2	14.7	13.9	13.4	13.5	14.2	14.9	15.0	14.6	14.3
10°	12.8	13.9	14.8	15.2	15.0	14.8	14.8	15.0	14.9	14.1	13.1	12.4
20°	10.8	12.3	13.9	15.2	15.7	15.8	15.7	15.3	14.4	12.9	11.2	10.3
30°	8.5	10.5	12.7	14.8	16.0	16.5	16.2	15.3	13.5	11.3	9.1	7.9
40°	6.0	8.3	11.0	13.9	15.9	16.7	16.3	14.8	12.2	9.3	6.7	5.4
50°	3.6	5.9	9.1	12.7	15.4	16.7	16.1	13.9	10.5	7.1	4.3	3.0

The parameter  $E_a$  is estimated as:

$$E_a = 0.35 \left( 1 + \frac{u_2}{160} \right) (e_w - e_a)$$

590 where  $u_2$  is the wind speed at 2m above ground in km/day;  $e_w$  is the saturation vapour pressure at mean air temperature in  
mm of mercury;  $e_a$  is the actual vapour pressure.

**Table A2:- Mean Monthly Values of Possible Sunshine Hours,  $N$**

North latitude	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
0°	12.1	12.1	12.1	12.1	12.1	12.1	12.1	12.1	12.1	12.1	12.1	12.1
10°	11.6	11.8	12.1	12.4	12.6	12.7	12.6	12.4	12.9	11.9	11.7	11.5
20°	11.1	11.5	12.0	12.6	13.1	13.3	13.2	12.8	12.3	11.7	11.2	10.9
30°	10.4	11.1	12.0	12.9	13.7	14.1	13.9	13.2	12.4	11.5	10.6	10.2
40°	9.6	10.7	11.9	13.2	14.4	15.0	14.7	13.8	12.5	11.2	10.0	9.4
50°	8.6	10.1	11.8	13.8	15.4	16.4	16.0	14.5	12.7	10.8	9.1	8.1



## Appendix B: Correlation analysis of catchment attributes

- 595 To explore the potential connections between various types of watershed attributes, we did correlation analysis using the Kendall rank correlation coefficient (Kendall 1938). Kendall rank correlation coefficient is a measure of rank correlation: the similarity of the sort order of the two sets of data. Kendall correlation will be high if the orderings of the observations of two variables are similar. Kendall correlation avoids the assumption of linear relationship and that the distribution should be normal and continuous (e.g., Pearson correlation coefficient; the results can be found in). When the relationship is not exactly linear, using Pearson correlation will miss out on information that Kendall could capture. Table B1, ~~which~~ shows the top five most relevant attributes for each attribute, ~~and the Fig. S1, the correlation matrix.~~ The analysis result shows that the correlations between variables are consistent in line with general understanding, justifying the rationality of the dataset, to name a few:
- 600 (1) Subsurface permeability and porosity are highly most correlated with geological attributes.
- ~~(2) LAI and NDVI have a high positive correlation (0.866).~~
- 605 ~~(3)(2) Root depth is are~~ most positively correlated with each other but most negatively correlated with the fraction of barren land cover types.
- ~~(3) Urban and built ups are most positively correlated with population density.~~
- (4) In China, the savanna is mainly distributed in the southern coastal areas, resulting in that it is most positively correlated with average rainfall (0.604)-mean precipitation.
- 610 (5) Sand is most positively correlated with the saturated hydraulic conductivity (0.86) while the clay is strongly negatively correlated (-0.763), and catchments with a lot of rainfall are less likely to have soil with high hydraulic conductivity (-0.647).
- ~~(6) High altitude catchments tend to have lower saturated water content (-0.705).~~

615 **Table B1:** The top five most relevant characteristics for each attribute (different soil layers for the same attribute are excluded, e.g., phihox\_sl2 is not included in the top five most relevant attributes of phihox\_sl1 though they are highly correlated)

Attribute	1 <sup>st</sup>	2 <sup>nd</sup>	3 <sup>rd</sup>	4 <sup>th</sup>	5 <sup>th</sup>
high_prec_fre q	<del>low_prec_dur</del> <u>root_dc</u> <u>pth_50(-0.58196)</u>	<del>root_depth_50(-</del> <u>grassland(0.438175)</u>	<del>root_depth_99(-</del> <u>0.436171)</u>	<del>barren(-</del> <u>som(0.39136)</u>	<del>pet_mean</del> <u>tk_satu_11(-</u> <u>0.26+133)</u>
high_prec_du r	<del>elev(theta_s_16(-</del> <u>0.544277)</u>	<del>theta_s_1615(-</del> <u>0.503234)</u>	<del>p_seasonality(pre</del> <u>mean(-0.49233)</u>	<del>theta_s_15(-</del> <u>elev(0.458211)</u>	<del>rhu_mean</del> <u>theta_s_14</u> <u>(-0.43+201)</u>
low_prec_fre q	<del>pre_mean(-</del> <u>0.88+766)</u>	<del>ssd_mean</del> <u>aridity(0.84</u> <u>+745)</u>	<del>phihox_sl7</del> <u>ssd_mea</u> <u>n(0.825652)</u>	<del>phihox_sl6(rhu_mea</del> <u>n(-0.818627)</u>	<del>phihox_sl5</del> <u>l7(0.81</u> <u>4588)</u>
low_prec_dur	<del>barren</del> <u>aridity(0.7287</u> <u>8)</u>	<del>rhu</del> <u>pre_mean(-</u> <u>0.723768)</u>	<del>exp</del> <u>ssd_mean(0.72</u> <u>+731)</u>	<del>ndv</del> <u>rhu_mean(-</u> <u>0.684709)</u>	<del>phihox_sl7</del> <u>root_d</u> <u>epth_99(0.66579)</u>

frac_snow_da ily	temgst_mean(- 0.951802)	gsttem_mean(- 0.949792)	ssd_meanlat(0.7775 75)	pre_meanevergreen _broadleaf tree(- 0.762512)	n_min(pre_mean(- 0.703436)
prs_mean p_seasonality	pre_mean(elev(- 0.904678)	rhu_meanlon(0.76555 2)	ssdrhu_mean(- (0.764432)	low_prec_freq(- 0.712)urban and b uilt- up land(0.427)	frac_snow_dailybarre n(-0.68341)
petpre_mean	cecsol_sl2aridity(- 0.66913)	cecsol_sl1low_prec dur(-0.634768)	cecsol_sl3low_prec freq(-0.628766)	gstssd_mean(- 0.622723)	bidfie_sl1rhu_mean (0.608712)
precvp_mean	p_seasonalityaridity(0 .904643)	low_prec_freqndvi_m ean(-0.881632)	ssdrhu_mean(- 0.858617)	rhu_ssd_mean(0.832 598)	phibox_sl7lai_dif(- 0.819593)
temgst_mean	gsttem_mean(0.9929 24)	frac_snow_daily(- 0.951802)	pre_mean(lat(- 0.747512)	ssd_mean(- 0.709)evergreen br oadleaf tree(0.50 7)	pet_meanp_season ality(0.681442)
prsrhu_mean	elevaridity(- 0.889751)	e_max(ssd_mean(- 0.707746)	lowpre_mean(0.70 7712)	e_min(low_prec_d ur(-0.707709)	rhu_mean(low_prec freq(-0.603627)
rhu pet_mean	ssd_meancecsol_sl2(- 0.887451)	pregst_mean(0.8324 42)	evp_meancecsol_sl 3(-0.823441)	ndvi_mean(cecsol_s l1(-0.813422)	low_prec_freqcecsol _sl4(-0.80342)
expssd_mean	ndvi_mean(- aridity(0.845753)	rhu_mean(- 0.823746)	ssd_meanlow_prec _dur(0.756731)	e_minpre_mean(- 0.731723)	ton(- low_prec_freq(0.7 3652)
win_mean	ssd_mean(0.581426 )	frac_snow_daily(wood y_savanna(- 0.571393)	tem_mean(- 0.52379)	gst_mean(- 0.507377)	low_prec_freq(mixed forest(-0.477363)
ssdtem_mean	rhu gst_mean(- (0.887924)	pre_meanfrac_snow daily(-0.858792)	low_prec_freqeverg reen_broadleaf t ree(0.84493)	frac_snow_dailypop dnsty(0.777475)	p_seasonalitylat(- 0.764474)
p_seasonality gst_mean	temrhu_mean(- 0.992421)	frac_snow_dailytem mean(-0.949397)	pregst_mean(- 0.743393)	n_min(- ssd_mean(0.69339 3)	lat(- low_prec_dur(0.69 3375)
aridity meability	pre_mean(- 0.408913)	sm(- low_prec_dur(0.403 78)	ssd_mean(0.399 753)	se(rhu_mean(- 0.323751)	bdtemlow_prec_fr eq(0.24745)

<u>slope</u>	<u>geol_porosity</u>	<u>su</u> (lat(-0.627374)	<u>pa</u> bdticm(-0.575348)	<u>phi</u> lox_sl1(win_me an(-0.46341)	<u>phi</u> lox_sl3mixed fo rest(0.454341)	<u>phi</u> lox_sl4evergree n_needleleaf tree(0.453327)
<u>iglon</u>	<u>snow_and_ice</u> (elev(-0.474585)	<u>tk</u> satu_15prs_mean(0.324552)	<u>tk</u> satu_13(cvp_mea n(-0.3485)	<u>tk</u> satu_14(barren(-0.306482)	<u>tk</u> satu_12ndvi_mean(0.27547)	
<u>elev</u>	<u>pa</u> geol_porosityprs_mea n(-0.575678)	<u>phi</u> lox_sl1lon(-0.344585)	<u>nd</u> _built-up_land(-0.302485)	<u>phi</u> lox_sl2pop_dnst y(-0.304481)	<u>phi</u> lox_sl4cropland(-0.297456)	
<u>se</u> lat	<u>geol_porosity</u> (-frac_snow_daily(0.362575)	<u>geol_permeability</u> (0.323)evergreen_broadle af_tree(-0.548)	<u>n_max</u> gst_mean(-0.347512)	<u>lat</u> tem_mean(-0.347474)	<u>n_min</u> (-low_prec_freq(0.346437)	
<u>su</u> pop	<u>geol_porosity</u> urban_a nd_built-up_land(0.627618)	<u>bd</u> tiemcropland(0.599519)	<u>cro</u> pland(aridity(-0.468511)	<u>phi</u> lox_sl1pre_mea n(0.44505)	<u>phi</u> lox_sl4rhu_mea n(0.439492)	
<u>pop</u> _dnsty	<u>sm</u> geol_permeability(-0.403)urban_and_bu ilt-up_land(0.639)	<u>su</u> aridity(-0.385538)	<u>cro</u> pland(-0.268533)	<u>bd</u> tiem(-pre_mean(0.233533)	<u>e_max</u> ssd_mean(-0.228521)	
<u>length</u>	<u>vi</u> deciduous_broadleaf treearea(0.244684)	<u>geol_porosity</u> form_fa ctor(-0.48398)	<u>lai_max</u> shape_fact or(0.465398)	<u>lai_diff</u> elongation_r atio(-0.459398)	<u>e_max</u> compactness coefficient(0.457363)	
<u>area</u>	<u>mt</u> geol_porosity(-length(0.442684)	<u>evergreen</u> _needleleaf treepop(0.32723)	<u>ore</u> dre_sl3pa(0.265194)	<u>ore</u> dre_sl4(circulato ry_ratio(-0.258187)	<u>bd</u> fie_sl5(-0.254)compactness coefficient(0.187)	
<u>ss</u> form_factor	<u>geol_permeability</u> (-0.408)elongation_r atio(1.0)	<u>su</u> (-shape_factor(-1.0)287)	<u>sm</u> (-circulatory_ratio(0.206435)	<u>geol_porosity</u> (0.2)CO mpactness coeffic ient(-0.435)	<u>tk</u> satu_16length(-0.456398)	
<u>pi</u> shape_factor	<u>deciduous</u> _broadleaf tree(0.299)elongation_ratio(-1.0)	<u>geol_porosity</u> (-0.208)form_factor(-1.0)	<u>e_max</u> (circulatory_ratio(-0.46435)	<u>lon</u> compactness_c oefficient(0.46435)	<u>e_min</u> length(0.46398)	
<u>va</u> compactnes s_coefficient	<u>geol_porosity</u> (-0.248)circulatory_ra tio(-1.0)	<u>high_pre</u> _dur(elongat ion_ratio(-0.494435)	<u>tem</u> _mean(-shape_factor(0.467435)	<u>gst</u> _meanform_fact or(-0.46435)	<u>su</u> (-length(0.46363)	

<u>circulatory_ratio</u>	<u>water_bodies_compactness_coefficient(-1.0)</u>	<u>permanent_wetland_elongation_ratio(0.379435)</u>	<u>shape_factor_root_depth_50(-0.464435)</u>	<u>theta_s_13_form_factor(0.448435)</u>	<u>theta_s_14_length(-0.447363)</u>
<u>elongation_ratio</u>	<u>theta_s_16_shape_factor(-1.0437)</u>	<u>theta_s_15_form_factor(1.0433)</u>	<u>elev(m)_circulatory_ratio(0.424435)</u>	<u>theta_s_14_compactness_coefficient(-0.444435)</u>	<u>pre_mean_length(-0.402398)</u>
<u>lai_diff</u>	<u>eesol_sl2_ndvi_mean(0.222808)</u>	<u>eesol_sl3_barren(-0.213642)</u>	<u>eesol_sl4_aridity(-0.212638)</u>	<u>eesol_sl4_pre_mean(0.211609)</u>	<u>eesol_sl5_woody_savanna(0.208607)</u>
<u>lai_max</u>	<u>snow_and_ice_ndvi_mean(0.206779)</u>	<u>theta_s_12_barren(-0.454614)</u>	<u>theta_s_13_aridity(-0.454613)</u>	<u>theta_s_11_woody_savanna(0.444612)</u>	<u>tk_satu_14(phiahox_sl2(-0.436602)</u>
<u>ndvi_mean</u>	<u>phiahox_sl1_lai_diff(0.214808)</u>	<u>phiahox_sl2_lai_max(0.207779)</u>	<u>phiahox_sl3_barren(-0.207677)</u>	<u>phiahox_sl4_cvp_mean(-0.205632)</u>	<u>phiahox_sl5_aridity(-0.202607)</u>
<u>root_depth_50_ev</u>	<u>tk_satu_13_grassland(-0.07485)</u>	<u>tk_satu_14_pet_mean(0.066232)</u>	<u>barren(0.064212)</u>	<u>high_prec_freq(tk_satu_12(-0.061196)</u>	<u>pdep(tk_satu_11(-0.061176)</u>
<u>root_depth_9_lai_diff</u>	<u>ndvi_mean(grassland)(-0.866339)</u>	<u>phiahox_sl4_barren(0.809337)</u>	<u>phiahox_sl2_cropland(-0.807336)</u>	<u>phiahox_sl5_pdep(-0.807284)</u>	<u>phiahox_sl6_lon(-0.807283)</u>
<u>lai_max_evergreen_needleleaf_tree</u>	<u>ndvi_mean(mixed_forest)(0.856572)</u>	<u>phiahox_sl4_woody_savanna(0.845481)</u>	<u>phiahox_sl5_sl7(-0.814416)</u>	<u>phiahox_sl6(-0.814411)</u>	<u>phiahox_sl2_sl5(-0.813409)</u>
<u>ndvi_mean_evergreen_broadleaf_tree</u>	<u>lai_diff(lat)(-0.866548)</u>	<u>lai_max(phiahox_sl7(-0.856538)</u>	<u>evp_mean(phiahox_sl16(-0.845529)</u>	<u>thu_mean(phiahox_sl5(-0.813522)</u>	<u>barren(pre_mean(0.772512)</u>
<u>root_depth_50_deciduous_needleleaf_tree</u>	<u>barren(cccsol_sl1)(0.56274)</u>	<u>low_prec_dur(bldfie_sl11(-0.626274)</u>	<u>grassland(cccsol_sl2)(0.537272)</u>	<u>orcdrc_sl2(ndvi_mean(-0.51327)</u>	<u>evp_mean(cccsol_sl3)(0.497262)</u>
<u>root_depth_90_deciduous_broadleaf_tree</u>	<u>barren(mixed_forest)(0.897604)</u>	<u>low_prec_dur(woody_savanna)(0.66568)</u>	<u>ndvi_mean(-0.628524)</u>	<u>evp_mean_lai_max(0.6045)</u>	<u>thu_mean(-lai_diff)(0.486497)</u>
<u>evergreen_needleleaf</u>	<u>woody_savanna_sl2(0.398713)</u>	<u>bldfie_sl4(-0.3914)deciduous_broadleaf_tree(0.604)</u>	<u>bldfie_sl5(-0.384)evergreen_needleleaf(0.604)</u>	<u>bldfie_sl3(phiahox_sl7(-0.372565)</u>	<u>bldfie_sl7(phiahox_sl6(-0.366563)</u>

<u>tree</u> <u>mixed_fore</u>			<u>eedleleaf_tree(0.</u>		
<u>st</u>			<u>572)</u>		
<u>evergreen</u> <u>broadleaf</u>	<u>pre_meandeciduous</u> <u>broadleaf_tree(0.50</u>	<u>lai_maxsavanna(0.48</u>	<u>phi_hox_sl7(-</u> <u>mixed_forest(0.4</u>	<u>lai_diff(tksatu_l4(-</u> <u>0.47+153)</u>	<u>phi_hox_sl6theta_s_l</u> <u>2(-0.47+142)</u>
<u>tree</u> <u>closed_shr</u> <u>ubland</u>	<u>4217)</u>	<u>3+16)</u>	<u>7+158)</u>		
<u>deciduous</u> <u>needleleaf</u>	<u>woody</u> <u>savanna_high_prec_d</u>	<u>eesol_sl2(rhu_mean(</u>	<u>eredre_sl2elev(0.22</u>	<u>petssd_mean(-</u>	<u>bl_dife_sl+prs_mean(</u>
<u>tree</u> <u>open_shru</u> <u>bland</u>	<u>ur(0.24+179)</u>	<u>-0.23+174)</u>	<u>6+17)</u>	<u>(0.21+517)</u>	<u>-0.21+165)</u>
<u>deciduous</u> <u>broadleaf</u>	<u>lai_maxmixed_forest</u>	<u>lai_diff(phi_hox_sl7(-</u>	<u>eesol_sl4(phi_hox</u>	<u>bl_dife_sl+phi_hox_sl</u>	<u>e_max(phi_hox_sl6(</u>
<u>tree</u> <u>woody_sav</u> <u>anna</u>	<u>(0.459+713)</u>	<u>0.452+628)</u>	<u>sl4(-0.433+628)</u>	<u>3(-0.413+627)</u>	<u>-0.36+627)</u>
<u>mixed</u> <u>forest</u> <u>savanna</u>	<u>eredre_sl+pre_mean(</u> <u>0.50+606)</u>	<u>lai_maxcropland_nat</u> <u>ural_vegetaion(0.47</u> <u>+605)</u>	<u>lai_diffwoody_sava</u> <u>nna(0.466+604)</u>	<u>phi_hox_sl6aridity(-</u> <u>0.462+602)</u>	<u>phi_hox_sl7ssd_mea</u> <u>n(-0.46+591)</u>
<u>closed</u> <u>shrubland</u> <u>grassl</u> <u>and</u>	<u>theta_s_sl+root_depth</u> <u>_50(-0.084+485)</u>	<u>grav(0.079)cropland</u> <u>natural_vegetaion(-</u> <u>0.363)</u>	<u>se(tem_mean(-</u> <u>0.075+344)</u>	<u>theta_s_sl2gst_mean(</u> <u>-0.072+344)</u>	<u>urban_and_built_up</u> <u>land(0.064)root_dep</u> <u>th_99(-0.339)</u>
<u>open</u> <u>shrubland</u> <u>perma</u> <u>nent_wetland</u>	<u>high_prec_durwater_b</u> <u>odies(0.155+469)</u>	<u>theta_s_sl6(-</u> <u>savanna(0.15+363)</u>	<u>rhu_mean(-</u> <u>0.149)urban_and</u> <u>built-</u> <u>up_land(0.347)</u>	<u>prspre_mean(-</u> <u>(0.147+343)</u>	<u>evp_meanpop(0.139</u> <u>343)</u>
<u>woody</u> <u>savanna</u> <u>croplan</u> <u>d</u>	<u>lai_maxurban_and_b</u> <u>uilt-</u> <u>up_land(0.633+546)</u>	<u>lai_diffpop_dnsty(0.6</u> <u>3+533)</u>	<u>phi_hox_sl7(-</u> <u>pop(0.592+519)</u>	<u>phi_hox_sl6elev(-</u> <u>0.594+56)</u>	<u>lon(phi_hox_sl5(-</u> <u>0.585+417)</u>
<u>savanna</u> <u>urban_a</u> <u>nd_built-</u> <u>up_land</u>	<u>pre_meanpop_dnsty(</u> <u>0.604+639)</u>	<u>phi_hox_sl7(-</u> <u>pop(0.556+18)</u>	<u>elaycropland(0.54</u> <u>7546)</u>	<u>phi_hox_sl6elev(-</u> <u>0.543+485)</u>	<u>phi_hox_sl5(-</u> <u>0.537)cropland_nat</u> <u>ural_vegetaion(0.</u> <u>428)</u>
<u>grassland</u> <u>cropla</u> <u>nd_natural_ve</u> <u>getaion</u>	<u>root_depth_50(-</u> <u>savanna(0.537+605)</u>	<u>temrhu_mean(-</u> <u>(0.496+546)</u>	<u>gst_meanaridity(-</u> <u>0.491+523)</u>	<u>frac_snow_daily(ssd</u> <u>mean(-0.469+52)</u>	<u>phi_hox_sl6pre_mea</u> <u>n(0.438+51)</u>

permanent wetlandsnow_and_ice	wbig(0.379431)	water bodiesbarren(0.349379)	p_seasonality(lon(-0.3373)	pre_meanelev(0.248369)	pop_dnsty(pdep(-0.23354)
eroplantbarren	su(ndvi_mean(-0.468677)	lon(lai dif(-0.412642)	e_min(lai_max(-0.412614)	e_maxaridity(0.412581)	evp_mean(elev(-0.388574)
urban-and-built-up landwater_bodies	pop_dnstypermanent_wetland(0.811469)	popWb(0.39939)	p_seasonalitycropland_natural_vegetation(0.28617)	tem_meanurban_and_built-up_land(0.261158)	elev(-0.244154)
geol_permeabilityeroplant/natural-vegetaion	ssd_meanSM(-0.458345)	savannaSU(0.381326)	rhu_mean(SS(-0.371316)	frac_snow_daily(bdticm(0.367228)	tem_meanpdep(0.364161)
geol_porositysnow-and-ice	tkstatu_15su(0.568455)	tkstatu_13(pa(-0.561417)	tkstatu_14(woody_savanna(-0.533323)	tkstatu_12(phiahox_sl3(0.506315)	tkstatu_14(phiahox_sl4(0.503314)
barrenig	root_depth_99snow_and_ice(0.897431)	root_depth_50elev(0.856194)	ndvi_mean(theta_sl2(-0.772185)	low_pree_dur(pdep(-0.728184)	evp_mean(theta_sl13(-0.698182)
water-bodiespa	wb(geol_porosity(-0.674417)	permanent_wetlandmt(0.3493)	root_depth_50(pi(0.192295)	root_depth_99(va(0.154271)	theta_sl13vi(0.153246)
sclength	area(geol_porosity(-0.849285)	eirculatory_ratio(lat(-0.491264)	elongation_ratio(bdticm(-0.45126)	form_factor(slope(0.436246)	compactness_coefficientmixed_forest(0.292231)
areaSU	lengthbdticm(0.84952)	popgeol_porosity(0.418455)	eirculatory_ratio(WOODY_SAVANNA(-0.255349)	eesol_sl1geol_permeability(0.142326)	bldfie_sl2(phiahox_sl7(0.138326)
form_factorSM	elongation_ratio(geol_permeability(-0.992345)	eirculatory_ratio(SU(-0.647283)	shape_factor(bdticm(-0.506228)	lengthcropland(-0.436199)	compactness_coefficientelev(0.194)
shape_factorVi	compactness_coefficientpa(0.786246)	elongation_ratio(pi(0.566203)	form_factor(va(0.506171)	eirculatory_ratio(geol_perosity(-0.372169)	lengthdeciduous_broadleaf_tree(0.266166)
compactness_coefficientmt	shape_factorpa(0.7863)	geol_porosityeirculatory_ratio(-0.594286)	elongation_ratio(pi(0.421199)	form_factor(deciduous_broadleaf_tree(0.187)	lengtharca(0.29218)

<u>circulatory_ratio_s</u>	<u>elongation_ratio</u> ( <u>gcol</u> <u>permeability</u> (- 0.654316)	<u>form_factor</u> ( <u>su</u> (- 0.64717)	<u>compactness</u> - <u>coeffic</u> <u>entbdticm</u> (- 0.594136)	<u>length</u> (- 0.491) <u>evergreen_ne</u> <u>edleleaf_tree</u> (0.10 6)	<u>shape_factor</u> ( <u>ksatu</u> 1 6(-0.372096)
<u>elongation_ratio</u> <u>i</u>	<u>form_factor</u> ( <u>pa</u> (0.9922 95)	<u>circulatory_ratio</u> ( <u>vi</u> (0.6 54203)	<u>shape_factor</u> (- <u>mt</u> (0.566199)	<u>length</u> ( <u>gcol</u> <u>porosit</u> <u>y</u> (-0.454183)	<u>compactness</u> - <u>coefficie</u> <u>nt</u> (-0.424) <u>va</u> (0.172)
<u>elev</u> ( <u>m</u> ) <u>va</u>	<u>prs_mean</u> (- <u>pa</u> (0.889271)	<u>e_min</u> ( <u>gcol</u> <u>porosity</u> ( -0.753219)	<u>lon</u> (- <u>vb</u> (0.75221)	<u>e_max</u> (- 0.752) <u>deciduous_n</u> <u>eedleleaf_tree</u> (0.1 86)	<u>theta_s</u> 14(- <u>pi</u> (0.7172)
<u>slope</u> ( <u>m/km</u> ) <u>wb</u>	<u>n_min</u> (- <u>water_bodies</u> (0.552 39)	<u>lat</u> (- <u>permanent_wetland</u> (0.554264)	<u>n_max</u> (- <u>bldfie_sl4</u> (0.5514 8)	<u>phihox_sl7</u> (- <u>bldfie_sl5</u> (0.49414 7)	<u>oredre_sl</u> 1 <u>urban_an</u> <u>d_built-</u> <u>up_land</u> (0.49138)
<u>n_min</u> <u>pb</u>	<u>lat</u> (1- <u>mt</u> (0.176)	<u>frac_snow_daily</u> ( <u>pa</u> (0.7 03132)	<u>gst_mean</u> ( <u>theta_s</u> 1 <u>5</u> (-0.693128)	<u>pre_mean</u> (- <u>area</u> (0.654127)	<u>tem_mean</u> (- <u>length</u> (0.648123)
<u>n_max</u> <u>vb</u>	<u>lat</u> (1- <u>va</u> (0.21)	<u>frac_snow_daily</u> ( <u>gcol</u> <u>porosity</u> (-0.704171)	<u>gst_mean</u> (- <u>vi</u> (0.692165)	<u>pre_mean</u> (- <u>cccsol_sl7</u> (0.6516 1)	<u>tem_mean</u> (- <u>cccsol_sl6</u> (0.64715 7)
<u>e_min</u> <u>nd</u>	<u>lon</u> (1- <u>barren</u> (0.154)	<u>elev</u> (- <u>aridity</u> (0.753146)	<u>evp</u> <u>pre_mean</u> (- 0.734144)	<u>prs_mean</u> ( <u>lai_dif</u> (- 0.707141)	<u>ndvi_mean</u> ( <u>snow_an</u> <u>d_ice</u> (0.694141)
<u>e_max</u> <u>py</u>	<u>lon</u> (1- <u>phihox_sl1</u> (- 0.237)	<u>elev</u> ( <u>phihox_sl2</u> (- 0.752233)	<u>evp_mean</u> ( <u>phihox_s</u> <u>13</u> (-0.729233)	<u>prs_mean</u> ( <u>phihox_sl</u> <u>4</u> (-0.70723)	<u>ndvi_mean</u> ( <u>woody_s</u> <u>avanna</u> (0.69227)
<u>pop</u> ( <u>people</u> ) <u>cv</u>	<u>area</u> ( <u>barren</u> (0.418036 )	<u>urban</u> — <u>and</u> — <u>built up</u> <u>land</u> (0.399) <u>orcdrc_sl5</u> (-0.035)	<u>tem_mean</u> ( <u>orcdrc_s</u> <u>14</u> (-0.348035)	<u>p_seasonality</u> ( <u>cccsol</u> <u>_sl3</u> (-0.347034)	<u>frac_snow_daily</u> ( <u>orcdrc</u> <u>c_sl7</u> (-0.304034)
<u>tk</u> <u>sat</u> <u>11</u> <u>pop</u> <u>dnst</u> ( <u>people/km</u> <u>2</u> )	<u>urban</u> — <u>and</u> — <u>built up</u> <u>land</u> (0.811) <u>grav</u> (- 0.346)	<u>p_seasonality</u> ( <u>som</u> (- 0.426344)	<u>tem_mean</u> ( <u>bldfie_sl</u> <u>3</u> (0.442298)	<u>gst_mean</u> ( <u>bldfie_sl1</u> ( 0.395295)	<u>frac_snow_daily</u> (- <u>bldfie_sl2</u> (0.39291 )
<u>tk</u> <u>sat</u> <u>12</u> <u>lon</u>	<u>e_max</u> (1- <u>som</u> (- 0.365)	<u>e_min</u> (1- <u>bldfie_sl3</u> (0. 326)	<u>elev</u> (- <u>bldfie_sl1</u> (0.7523 26)	<u>evp_mean</u> (- <u>bldfie_sl2</u> (0.73323 )	<u>prs_mean</u> ( <u>grav</u> (- 0.707308)
<u>lat</u> <u>tk</u> <u>sat</u> <u>13</u>	<u>n_min</u> (1- <u>som</u> (-0.344)	<u>n_max</u> (1- <u>bldfie_sl2</u> (0 .328)	<u>frac_snow_daily</u> ( <u>bldf</u> <u>ie_sl1</u> (0.702325)	<u>gst_mean</u> (- <u>bldfie_sl3</u> (0.69332 4)	<u>pre_mean</u> (- <u>bldfie_sl4</u> (0.65430 8)

tk satu_114	snow and sand(0.5033 98)	silt(-0.465397)	sand(- bldfie_sl1(0.3663 88)	sandbldfie_sl3(0.36 2384)	log_k_s_15bldfie_sl 4(0.327358)
tk satu_125	snow and sand(0.5063 86)	silt(- bldfie_sl2(0.49376)	sand(som(- 0.406369)	som(- bldfie_sl4(0.36536 4)	log_k_s_15bldfie_sl 1(0.364358)
tk satu_136	snow and sand(0.5643 66)	silt(- som(-0.489362)	sand(bdticm(0.4093 6)	ndvi_mean(- bldfie_sl2(0.36834 3)	clay(- bldfie_sl7(0.33433 8)
tk satu_14log_k_s_11	snow and sand(0.53371)	silt(- clay(-0.4959)	sand(savanna(- 0.465441)	ndvi_mean(- silt(- 0.455436)	log_k_s_15(rhu_mea n(-0.414423)
tk satu_15log_k_s_12	snow and sand(0.568709)	silt(- clay(-0.402578)	ndvi_mean(savanna (-0.375452)	sand(phiaox_sl7(0.3 48438)	clay(- silt(- 0.326433)
tk satu_16log_k_s_13	snow and sand(0.449682)	bdticm(- clay(- 0.403592)	log_k_s_16(savanna (-0.384448)	su(phiaox_sl7(0.38 442)	low_pre_freep(phiaox _sl6(0.36435)
log_k_s_114	sand(0.858612)	clay(-0.733603)	pre_mean(savanna(- 0.55349)	phiaox_sl7(pre_mea n(-0.554489)	phiaox_sl6sl7(0.54 6485)
log_k_s_125	sand(- clay(-0.86561)	clay(- sand(0.729555)	phiaox_sl7(0.575 506)	phiaox_sl6(savanna( -0.569501)	pre_mean(- phiaox_sl6(0.5685 01)
log_k_s_136	sand(- clay(-0.859563)	clay(- pre_mean(- 0.728555)	pre_mean(- aridity(0.571548)	phiaox_sl7(0.5715 34)	phiaox_sl6(0.5655 32)
log_k_theta_s_14_11	sand(- grav(-0.82582)	clay(- (-0.752325)	pre_meansand(- 0.647315)	phiaox_sl7(elev(- 0.636314)	phiaox_sl6pdep(0.63 311)
log_k_theta_s_15_12	sand(- grav(- 0.773585)	clay(- pdep(0.714377)	phiaox_sl7(elev(- 0.654366)	clay(phiaox_sl6(0. 64935)	phiaox_sl5(sand(- 0.646326)
log_k_theta_s_16_13	grav(- sand(0.688522)	clay(- pdep(0.68742)	phiaox_sl7(elev(- 0.665414)	phiaox_sl6prs_mean (0.662365)	pre_mean(- clay(0.662359)
theta_s_114	grav(-0.705515)	elev(- pdep(0.422463)	rhu_mean(elev(- 0.407412)	clayprs_mean(0.40 1349)	pdep(lon(0.4328)
theta_s_125	grav(-0.713433)	elev(-0.505401)	pdep(0.475376)	e_min(sand(- 0.442349)	low_rhu_mean(0.441 331)



<u>theta_s_13l6</u>	<u>grav(-</u> <u>0.662)evergreen_bro</u> <u>adleaf_tree(0.372)</u>	<u>elevgrav(-0.638357)</u>	<u>pr_mean(elev(-</u> <u>0.554344)</u>	<u>pdep(sand(-</u> <u>0.52343)</u>	<u>e_min(em_mean(0.</u> <u>516337)</u>
<u>theta_s_14orcdr</u> <u>c_sl7</u>	<u>elevbldfie_sl4(-</u> <u>0.7581)</u>	<u>gravbldfie_sl5(-</u> <u>0.663572)</u>	<u>pr_mean(bldfie_sl</u> <u>6(-0.594548)</u>	<u>pdep(bldfie_sl3(-</u> <u>0.574535)</u>	<u>e_min(bldfie_sl7(-</u> <u>0.51523)</u>
<u>theta_s_15orcdr</u> <u>c_sl3</u>	<u>elevbldfie_sl3(-</u> <u>0.656738)</u>	<u>gravbldfie_sl2(-</u> <u>0.584728)</u>	<u>pr_mean(bldfie_sl</u> <u>1(-0.536701)</u>	<u>pdep(bldfie_sl4(-</u> <u>0.504691)</u>	<u>flu_mean(bldfie_sl</u> <u>5(-0.467621)</u>
<u>theta_s_16orcdr</u> <u>c_sl4</u>	<u>elevbldfie_sl3(-</u> <u>0.637702)</u>	<u>pr_mean(bldfie_sl2(</u> <u>-0.525682)</u>	<u>gravbldfie_sl4(-</u> <u>0.513676)</u>	<u>bldfie_sl1high pr</u> <u>ee_dur(-0.503657)</u>	<u>flu_mean(bldfie_sl</u> <u>5(-0.475614)</u>
<u>orcdr_c17s15</u>	<u>eesol_sl2(bldfie_sl4</u> <u>(-0.758641)</u>	<u>bldfie_sl2sl3(-</u> <u>0.745636)</u>	<u>bldfie_sl4sl2(-</u> <u>0.744611)</u>	<u>bldfie_sl5sl5(-</u> <u>0.7376)</u>	<u>eesol_sl3(bldfie_sl</u> <u>1(-0.735592)</u>
<u>orcdr_c13s16</u>	<u>bldfie_sl2sl4(-</u> <u>0.876584)</u>	<u>bldfie_sl4sl5(-</u> <u>0.875567)</u>	<u>bldfie_sl3sl6(-</u> <u>0.874556)</u>	<u>bldfie_sl5sl3(-</u> <u>0.849552)</u>	<u>bldfie_sl1sl7(-</u> <u>0.848534)</u>
<u>orcdr_c14s12</u>	<u>bldfie_sl4sl2(-</u> <u>0.823787)</u>	<u>bldfie_sl2sl1(-</u> <u>0.809769)</u>	<u>bldfie_sl3(-</u> <u>0.803749)</u>	<u>bldfie_sl5sl4(-</u> <u>0.80368)</u>	<u>bldfie_eesol_sl1(-</u> <u>(0.787629)</u>
<u>orcdr_c15s11</u>	<u>bldfie_sl4phihox_sl2</u> <u>(-0.759599)</u>	<u>bldfie_sl2phihox_sl3</u> <u>(-0.754594)</u>	<u>bldfie_sl5phihox_s</u> <u>14(-0.745591)</u>	<u>bldfie_sl1phihox_sl</u> <u>5(-0.745586)</u>	<u>bldfie_sl3phihox_sl</u> <u>6(-0.731585)</u>
<u>oredre_sl6phihox</u> <u>x_sl7</u>	<u>woody_savanna(-</u> <u>eesol_sl2(0.733628</u> <u>)</u>	<u>bldfie_sl4pre_mean(-</u> <u>0.733598)</u>	<u>bldfie_sl2(-</u> <u>aridity(0.728592)</u>	<u>bldfie_sl1(-</u> <u>low_prec_freq(0.7</u> <u>25588)</u>	<u>orcdr_sl1bldfie_sl</u> <u>15(-0.721583)</u>
<u>phihox_sl6oredre</u> <u>dre_sl2</u>	<u>bldfie_sl2woody_sav</u> <u>anna(-0.917627)</u>	<u>bldfie_sl4pre_mean(-</u> <u>0.908594)</u>	<u>bldfie_sl3(-</u> <u>aridity(0.86459)</u>	<u>eesol_sl1(lai_max(</u> <u>-0.854587)</u>	<u>bldfie_sl4orcdr_sl1</u> <u>(-0.854585)</u>
<u>phihox_sl5oredre</u> <u>dre_sl4</u>	<u>phihox_sl2woody_sav</u> <u>anna(-0.826626)</u>	<u>phihox_sl4lai_max(-</u> <u>0.824593)</u>	<u>phihox_sl3pre_mean</u> <u>n(-0.822592)</u>	<u>phihox_sl4(-</u> <u>aridity(0.819589)</u>	<u>phihox_sl5orcdr_sl</u> <u>1(-0.813586)</u>
<u>phihox_sl7sl4</u>	<u>low_prec_freq(woody</u> <u>savanna(-</u> <u>0.825628)</u>	<u>pre_meanlai_max(-</u> <u>0.819599)</u>	<u>lai_maxorcdr_sl1</u> <u>(-0.806591)</u>	<u>oredre_sl1lai_dif(-</u> <u>0.804578)</u>	<u>lai_difpre_mean(-</u> <u>0.799576)</u>
<u>phihox_sl6sl3</u>	<u>low_prec_freq(woody</u> <u>savanna(-</u> <u>0.818627)</u>	<u>lai_max(-0.814595)</u>	<u>pre_meanorcdr_sl</u> <u>1(-0.81594)</u>	<u>oredre_sl1lai_dif(-</u> <u>0.807576)</u>	<u>lai_difpre_mean(-</u> <u>0.807568)</u>
<u>phihox_sl5sl2</u>	<u>lai_maxwoody_sav</u> <u>anna(-0.814627)</u>	<u>low_prec_freq(lai_max</u> <u>x(-0.814602)</u>	<u>orcdr_sl1(-</u> <u>0.813599)</u>	<u>lai_dif(-0.807583)</u>	<u>pre_mean(-</u> <u>low_prec_freq(0.8</u> <u>04569)</u>

phihox_s4s11	oredre_s11woody_sa vanna(-0.819601)	lai_max(-0.815586)	lai_diffordrc_sl1(- 0.809584)	low_prec_freq(lai_di f(-0.804565)	pre_mean(- bldfie_sl2(0.78155 )
phihox_s13bldfi e_sl7	ordrc_s15(- 0.822547)	lai_maxordrc_sl4(- 0.813546)	lai_diffordrc_sl3(- 0.806543)	ordrc_sl6(- low_prec_freq(0.7 99534)	pre_meanordrc_sl7 (-0.772523)
phihox_s12bldfi e_sl6	ordrc_s15(- 0.826559)	lai_maxordrc_sl6(- 0.813556)	lai_diffordrc_sl4(- 0.807553)	low_prec_freq(ordrc c_sl7(-0.798548)	pre_meanordrc_sl3 (-0.767547)
bldfie_sl5phihox s_s1	ordrc_s13(- 0.824621)	lai_maxordrc_sl4(- 0.804614)	lai_diffordrc_sl5(- 0.7986)	low_prec_freq(ordrc c_sl2(-0.78597)	pre_meanordrc_sl7 (-0.741572)
bldfie_s17s14	ordrc_sl3(- 0.775691)	ordrc_s14s12(- 0.74768)	ordrc_s15s14(- 0.698676)	ordrc_s12s15(- 0.698641)	ordrc_sl6(- 0.671584)
bldfie_s16s11	ordrc_s13s12(- 0.776769)	ordrc_s14s13(- 0.748701)	oredre_s15cecsol_sl 1(-0.701686)	ordrc_s12s14(- 0.694657)	oredre_s16som(- 0.677606)
bldfie_s15s13	ordrc_s13s12(- 0.849749)	ordrc_s12s13(- 0.81738)	ordrc_sl4(- 0.803702)	ordrc_sl5(- 0.745636)	oredre_s17som(- 0.728633)
bldfie_s14s12	ordrc_s13s12(- 0.875787)	ordrc_s12s13(- 0.854728)	ordrc_sl4(- 0.823682)	cecsol_sl1(- 0.763671)	oredre_s15som(- 0.759651)
bldfiececsol_sl 1	oredre_s12bldfie_sl1( -0.908686)	cecsol_s11bldfie_sl2( -0.891671)	ordrc_s13(- sl2(0.848629)	cecsol_s12bldfie_sl3 (-0.828598)	ordrc_s14(- sl3(0.787579)
cecsol_sl2bldfi e_s13	oredre_s13bldfie_sl1( -0.874579)	oredrebldfie_sl2(- 0.861566)	ordrc_s14(- sl2(0.803553)	cecsol_s11(- ordrc_sl3(0.7955 23)	sombldfie_sl3(- 0.787515)
bldfie_s12cecsol _sl5	oredre_s12bldfie_sl1( -0.917445)	oredre_s13bldfie_sl2( -0.876429)	cecsol_s11(- ordrc_sl2(0.874 12)	ordrc_s14(- sl3(0.809393)	sompct_mean(- 0.808392)
cecsol_s11s14	bldfie_sl1(- 0.891472)	bldfie_sl2(- 0.87459)	ordrc_sl2(0.8544 47)	bldfieordrc_sl3(- (0.79543)	ordrc_s13s15(0.781 424)
cecsol_s12s13	bldfie_sl1(- 0.828532)	oredrebldfie_sl2(- 0.82252)	bldfieordrc_sl2(- (0.798508)	ordrc_s17s13(0.758 49)	ordrc_s13s14(0.746 478)
cecsol_s15s17	bldfie_sl1(- 0.681413)	oredrebldfie_sl2(- 0.664396)	ordrc_s17s12(0.64 938)	bldfie_s12pet_mean( -0.645374)	ordrc_s16s13(0.636 362)
cecsol_s14s16	bldfie_sl1(- 0.72409)	oredrebldfie_sl2(- 0.717393)	ordrc_s17s12(0.69 3378)	bldfie_s12pet_mean( -0.692373)	ordrc_s16s13(0.679 36)

<u>eesol_sl3bdtic</u> <u>m</u>	<u>bldfie_sl1(-</u> <u>su(0.78452)</u>	<u>oredre_sl2(woody sa</u> <u>vanna(-0.776412)</u>	<u>bldfie_sl2(-</u> <u>low_prec_freq(0.</u> <u>76382)</u>	<u>oredre_phihox_sl7(0</u> <u>.735378)</u>	<u>oredre_sl3(mixed_f</u> <u>orest(-0.733374)</u>
<u>eesol_sl7pdep</u>	<u>bldfie_sl1(-</u> <u>theta_s_l4(0.661463</u> <u>)</u>	<u>oredre_sl7(elev(-</u> <u>0.654436)</u>	<u>oredre_sl2(grav(-</u> <u>0.642424)</u>	<u>oredre_sl6(theta_s_l</u> <u>3(0.6442)</u>	<u>oredre_sl5(lon(0.6194</u> <u>)</u>
<u>eesol_sl6por</u>	<u>bldfie_sl1(-</u> <u>som(0.648363)</u>	<u>oredre_sl2(bldfie_sl1</u> <u>(-0.637335)</u>	<u>oredre_sl7(phihox</u> <u>sl1(-0.62329)</u>	<u>oredre_sl6(phihox_s</u> <u>l3(-0.62328)</u>	<u>bldfie_phihox_sl2(-</u> <u>0.61328)</u>
<u>claybdtiem</u>	<u>su(sand(-0.59967)</u>	<u>low_prec_freq(log_k</u> <u>s_l4(-0.463603)</u>	<u>log_k_s_l6(l3(-</u> <u>0.439592)</u>	<u>phihox_sl2(log_k_s</u> <u>l1(-0.43759)</u>	<u>phihox_sl7(log_k_s</u> <u>l2(-0.436578)</u>
<u>sandpdep</u>	<u>elev(-</u> <u>log_k_s_l1(0.66271</u> <u>)</u>	<u>theta_log_k_s_l4(0.</u> <u>571709)</u>	<u>e_min_log_k_s_l3(</u> <u>0.566682)</u>	<u>lon(clay(-0.56567)</u> <u>.564612)</u>	<u>e_max_log_k_s_l4(0</u> <u>.564612)</u>
<u>siltpor</u>	<u>silt(sand(-0.573)</u>	<u>log_k_s_l1(-</u> <u>clay(0.366436)</u>	<u>ksatu_log_k_s_l2(-</u> <u>0.317433)</u>	<u>som(log_k_s_l3(-</u> <u>0.3144)</u>	<u>ksatu_l1_log_k_s_l4</u> <u>(-0.309316)</u>
<u>claygrav</u>	<u>pre_mean(theta_s_l2(</u> <u>-0.763585)</u>	<u>log_k_theta_s_l4(1(-</u> <u>0.752582)</u>	<u>log_k_theta_s_l1(3(</u> <u>-0.733522)</u>	<u>log_k_theta_s_l2(4(-</u> <u>0.729515)</u>	<u>log_k_theta_s_l3(5(-</u> <u>0.728433)</u>
<u>sandsom</u>	<u>log_k_s_l2(bldfie_sl2</u> <u>(-0.86651)</u>	<u>log_k_s_l3(bldfie_sl3</u> <u>(-0.859633)</u>	<u>log_k_s_l4(bldfie_s</u> <u>l1(-0.858606)</u>	<u>log_k_s_l4(orcdr_c_sl</u> <u>2(0.82599)</u>	<u>log_k_s_l5(orcdr_c_sl</u> <u>3(0.773576)</u>
<u>silt_high_prec_f</u> <u>req</u>	<u>por(root_depth_50(-</u> <u>0.573196)</u>	<u>sand(-</u> <u>grassland(0.558175)</u>	<u>log_k_s_l3(root de</u> <u>pth_99(-</u> <u>0.557171)</u>	<u>log_k_s_l2(-</u> <u>som(0.547136)</u>	<u>log_k_stksatu_l1(-</u> <u>0.545133)</u>
<u>high_prec_du</u> <u>rgrav</u>	<u>theta_s_l2(6(-</u> <u>0.713277)</u>	<u>theta_s_l1(5(-</u> <u>0.705234)</u>	<u>theta_s_l4(-</u> <u>p_seasonality(0.6</u> <u>63233)</u>	<u>theta_s_l3(-</u> <u>elev(0.662211)</u>	<u>theta_s_l5(4(-</u> <u>0.584201)</u>
<u>low_prec_fre</u> <u>qsom</u>	<u>bldfie_sl2(pre_mean(-</u> <u>0.808766)</u>	<u>bldfie_sl3(-</u> <u>aridity(0.787745)</u>	<u>bldfie_sl1(-</u> <u>ssd_mean(0.7596</u> <u>52)</u>	<u>bldfie_sl4(rhu_mean</u> <u>(-0.747627)</u>	<u>oredre_sl2(phihox_sl</u> <u>7(0.74588)</u>

## **Appendix C: Data sources and data processing**

The program to generate the data set is mainly written in Python. The rasterio<sup>13</sup> library is used to extract from the raster for the given basin boundary, reproject and merge rasters; The shapely<sup>14</sup> library is used to calculate the geometry; The pyproj<sup>15</sup> library is used for coordinate system conversions; The richdem<sup>16</sup> library is used to calculate slope; The netCDF4<sup>17</sup> and xarray<sup>18</sup> library is used to read the netCDF files; The pyshp<sup>19</sup> library is used to handle shapefiles; The gdal<sup>20</sup> command-line programs are used for data format conversions; The Python multiprocessing<sup>21</sup> library is used for multi-threaded data processing such as the calculation of meteorological time series; The interpolation program is written based on SciPy and NumPy. In addition, the calculation of the catchment boundary uses ArcPy<sup>22</sup>. However, ArcPy is not open sourced. The SURF\_CLI\_CHN\_MUL\_DAY dataset can be downloaded from [https://data.cma.cn/data/cdcdetail/dataCode/SURF\\_CLI\\_CHN\\_MUL\\_DAY.html](https://data.cma.cn/data/cdcdetail/dataCode/SURF_CLI_CHN_MUL_DAY.html). It is freely available to global researchers but registration is required. The GDBD dataset can be downloaded at [https://www.cger.nies.go.jp/db/gdbd/gdbd\\_index\\_e.html](https://www.cger.nies.go.jp/db/gdbd/gdbd_index_e.html). ASTER GDEM dataset can be downloaded at: <https://asterweb.jpl.nasa.gov/gdem.asp>. GLHYMPS dataset can be downloaded at: <https://dataverse.scholarsportal.info/dataset.xhtml?persistentId=doi:10.5683/SP2/DLGXYO>; MODIS MCD12Q1 can be obtained from: <https://lpdaac.usgs.gov/products/mcd12q1v006/>; MODIS MCD15A3 can be obtained from: <https://lpdaac.usgs.gov/products/mcd15a3hv006/>; Soil hydraulic and thermal properties can be downloaded after registration: <http://globalchange.bnu.edu.cn/research/soil5.jsp>; Soil properties data can be downloaded after registration: <http://globalchange.bnu.edu.cn/research/soil2>; SoilGrids250m data download links: <https://files.isric.org/soilgrids/former/2017-03-10/data/> with a list of descriptions: [https://github.com/ISRICWorldSoil/SoilGrids250m/blob/master/grids/models/META\\_GEOTIFF\\_1B.csv](https://github.com/ISRICWorldSoil/SoilGrids250m/blob/master/grids/models/META_GEOTIFF_1B.csv).

## **Appendix D: Basin boundaries**

This section briefly introduces how the basin boundaries are derived. The basin boundaries data used in this research are obtained from the GBDB (Masutomi, Inui et al. 2009) dataset. The GBDB dataset first distinguishing sinks caused by DEM

<sup>13</sup> <https://rasterio.readthedocs.io/en/latest/>

<sup>14</sup> <https://shapely.readthedocs.io/en/stable/manual.html>

<sup>15</sup> <https://pyproj4.github.io/pyproj/stable/>

<sup>16</sup> <https://richdem.readthedocs.io/en/latest/>

<sup>17</sup> <https://unidata.github.io/netcdf4-python/>

<sup>18</sup> <http://xarray.pydata.org/en/stable/>

<sup>19</sup> <https://pypi.org/project/pyshp/>

<sup>20</sup> <https://gdal.org/api/python.html>

<sup>21</sup> <https://docs.python.org/3/library/multiprocessing.html>

<sup>22</sup> <https://pro.arcgis.com/zh-cn/pro-app/latest/arcpy/get-started/what-is-arcpy-.htm>

640 errors, then the stream burning (Maidment 1996), and ridge fencing methods are used to modify the seeded DEM, then basin boundaries are produced with standardized procedures (Jenson, Domingue et al. 1988, Maidment and Morehouse 2002). Then the gauging station data from the GRDC (Center 2005) dataset is used to calibrate the derived basin boundaries. The derived basin areas were compared with the observed basin areas, and they showed a high degree of consistency with the observed basin data.

**Appendix E: Guidelines for generating basin attributes for any basin**

645 The published code<sup>23</sup> supports the automation of the calculation of the attributes for any given river basin and the generation of statistics files. In general, the user only needs to prepare the source data and ensure that the code environment is installed correctly, and then the user can run the code to calculate all attributes for the given river basin. The following describes the steps to generate data for any given watershed.

**Prepare source data**

650 In this step, the user needs to download the source data and place it in the corresponding location (Table D1). The code supports the calculation of meteorological time series based on the SURF\_CLI\_CHN\_MUL\_DAY data set. If the basin the user need to calculate is not in China, then the user needs to format the collected meteorological time series into the same format as the time series generated by the code. A sample file is available in the GitHub library.

655 **Table D1: Instructions for preparing data sources**

<u>Data source</u>	<u>Download link</u>	<u>Example</u>	<u>Note</u>
<u>ASTER</u>	<u><a href="https://search.earthdata.nasa.gov/search/">https://search.earthdata.nasa.gov/search/</a></u>	<u><a href="#">./data/dems/ *.tif</a></u>	
<u>GDEM</u>	<u><a href="https://www.jspacesystems.or.jp/ersdac/GDEM/E/">https://www.jspacesystems.or.jp/ersdac/GDEM/E/</a></u>		
<u>GLHYMP</u>	<u><a href="https://dataverse.scholarsportal.info/dataset.xhtml?persistentId=doi:10.5683/SP2/DLGXYO">https://dataverse.scholarsportal.info/dataset.xhtml?persistentId=doi:10.5683/SP2/DLGXYO</a></u>	<u><a href="#">./data/processed permeability.tif</a></u>	
<u>S</u>	<u><a href="#">fo/dataset.xhtml?persistentId=doi:10.5683/SP2/DLGXYO</a></u> (using <u><a href="#">source data requires merging multiple small pieces to a single TIFF</a></u> )	<u><a href="#">./data/processed porosity.tif</a></u>	

<sup>23</sup> <https://github.com/haozhen315/CCAM-China-Catchment-Attributes-and-Meteorology-dataset>

---

	<a href="https://1drv.ms/u/s!AqzR0fLyn9KKspF6HAAuXU9Twwkz1Q?e=QCPFAM">https://1drv.ms/u/s!AqzR0fLyn9KKspF6HAAuXU9Twwkz1Q?e=QCPFAM</a> (our processed file)		
	<a href="https://1drv.ms/u/s!AqzR0fLyn9KKspF70EPmDubS5V2qTQ?e=Rbybwa">https://1drv.ms/u/s!AqzR0fLyn9KKspF70EPmDubS5V2qTQ?e=Rbybwa</a> (our processed file)		
<u>GLiM</u>	<a href="https://csdms.colorado.edu/wiki/Data:GLiM">https://csdms.colorado.edu/wiki/Data:GLiM</a>	<a href="#">./data/processed_glim.py</a>	
	<a href="https://1drv.ms/u/s!AqzR0fLyn9KKspF5Vktb-zlmd_Ctxg?e=G6fOuh">https://1drv.ms/u/s!AqzR0fLyn9KKspF5Vktb-zlmd_Ctxg?e=G6fOuh</a> (our processed file)		
<u>MCD12Q1</u>	<a href="https://lpdaac.usgs.gov/products/mcd12q1v006/">https://lpdaac.usgs.gov/products/mcd12q1v006/</a>	<a href="#">./data/processed_igbp.tif</a>	
	<a href="https://1drv.ms/u/s!AqzR0fLyn9KKspF4xxbe0xM7qJNzka?e=vvFcFj">https://1drv.ms/u/s!AqzR0fLyn9KKspF4xxbe0xM7qJNzka?e=vvFcFj</a> (our processed file)		
<u>MCD15A3</u>	<a href="https://lpdaac.usgs.gov/products/mcd15a3hv006/">https://lpdaac.usgs.gov/products/mcd15a3hv006/</a>	<a href="#">./data/MCD15A3/MCD15A3H.A2002185.h22v04.006.2015149102803.hdf</a>	
<u>MOD13Q1</u>	<a href="https://lpdaac.usgs.gov/products/mod13q1v006/">https://lpdaac.usgs.gov/products/mod13q1v006/</a>	<a href="#">./data/MOD13Q1/MOD13Q1.A2002186.h22v04.006.2015149102803.hdf</a>	
<u>Soil</u>	<a href="http://globalchange.bnu.edu.cn/research/soil5.jsp">http://globalchange.bnu.edu.cn/research/soil5.jsp</a>	<a href="#">./data/soil_souce_data/binary/log_k_s_ll</a>	
<u>Soil</u>	<a href="https://files.isric.org/soilgrids/former/2017-03-10/data/">https://files.isric.org/soilgrids/former/2017-03-10/data/</a>	<a href="#">./data/soil_souce_data/tif/BD_TICM_M_250m_ll.tif</a>	<u>Description:</u> <a href="https://github.com/ISRICWorldSoil/SoilGrids250m/blob/master/grids/metadata/META_GEOTIFF_1B.csv">https://github.com/ISRICWorldSoil/SoilGrids250m/blob/master/grids/metadata/META_GEOTIFF_1B.csv</a>
<u>Soil</u>	<a href="http://globalchange.bnu.edu.cn/research/soil2">http://globalchange.bnu.edu.cn/research/soil2</a>	<a href="#">./data/soil_souce_data/tif/SANC</a>	

---

<u>SURF_CLI_CHN_MUL_DAY</u>	<a href="https://data.cma.cn/data/cdcdetail/dataCode/SURF_CLI_CHN_MUL_DAY.html">https://data.cma.cn/data/cdcdetail/dataCode/SURF_CLI_CHN_MUL_DAY.html</a>	<a href="#">./data/SURF_CLI_CHN_MUL_DAY-EVP-13240-195101.TXT</a>	If basin boundary is outside China, format and place the collected time series data in <a href="#">./output/catchment_meteorological</a>
Root depth	<a href="https://github.com/haozhen315/CAM-China-Catchment-Attributes-and-Meteorology-dataset/blob/main/data/root_depth_calculated.txt">https://github.com/haozhen315/CAM-China-Catchment-Attributes-and-Meteorology-dataset/blob/main/data/root_depth_calculated.txt</a>	<a href="#">./data/root_depth_calculated.txt</a>	Calculated root depth of each land type according to (Zeng 2001).
GLiM name mapping	<a href="https://github.com/haozhen315/CAM-China-Catchment-Attributes-and-Meteorology-dataset/blob/main/data/glim_name_short_long.txt">https://github.com/haozhen315/CAM-China-Catchment-Attributes-and-Meteorology-dataset/blob/main/data/glim_name_short_long.txt</a>	<a href="#">./data/glim_cate_number_mapping.csv</a> <a href="#">./data/glim_name_short_long.txt</a>	These files are used for name conversions in the program.
GDBD	<a href="https://www.cger.nies.go.jp/db/gdbd/gdbd_index_e.html">https://www.cger.nies.go.jp/db/gdbd/gdbd_index_e.html</a>	<a href="#">./data/river_network/as_streams_wgs.shp</a>	River network shapefiles are used to determine river basin shape factors. The source data need to be reprojected to EPSG:4326 (using ArcMap or QGIS) to successfully run the code. Note that files in different regions have different names.

### **Run the code**

When all the data is ready, the user can run the code [calculate\\_all\\_attributes.py](#) to calculate all attributes or run separate scripts (e.g., [soil.py](#)) to calculate indicators for specific categories. The result will appear in the output folder.

## 660 Financial support

This research has been supported by the National Key Research and Development Program (2018YFC0407901, 2018YFC0407905), the National Natural Science Fund of China (51779100), and the Central Public-interest Scientific Institution Basal Research Fund (HKY-JBYW-2020-21, HKY-JBYW-2020-07).

## References

- 665 [Abu Mostafa, Y. S., M. Magdon Ismail and H. T. Lin \(2012\). \*Learning from data\*, AMLBook New York, NY, USA.](#)  
[Abrams, M., R. Crippen and H. J. R. S. Fujisada \(2020\). "ASTER global digital elevation model \(GDEM\) and ASTER global water body dataset \(ASTWBD\)." \*Hydrology and Earth System Sciences \(HESS\)\* \*\*24\*\*\(7\): 1156.](#)
- Addor, N., H. X. Do, C. Alvarez-Garreton, G. Coxon, K. Fowler and P. A. J. H. S. J. Mendoza (2020). "Large-sample hydrology: recent progress, guidelines for new datasets and grand challenges." *Hydrology and Earth System Sciences (HESS)* **24**(5): 712-725.
- 670 Addor, N., A. J. Newman, N. Mizukami and M. P. Clark (2017). "The CAMELS data set: catchment attributes and meteorology for large-sample studies." *Hydrology and Earth System Sciences (HESS)* **21**(10): 5293-5313.
- Alvarez-Garreton, C., P. A. Mendoza, J. P. Boisier, N. Addor, M. Galleguillos, M. Zambrano-Bigiarini, A. Lara, G. Cortes, R. Garreaud and J. McPhee (2018). "The CAMELS-CL dataset: catchment attributes and meteorology for large sample studies-Chile dataset." *Hydrology and Earth System Sciences* **22**(11): 5817-5846.
- 675 Belward, A. S., J. E. Estes and K. D. Kline (1999). "The IGBP-DIS global 1-km land-cover data set DISCover: A project overview." *Photogrammetric Engineering and Remote Sensing* **65**(9): 1013-1020.
- Berghuijs, W. R., E. E. Aalbers, J. R. Larsen, R. Trancoso and R. A. Woods (2017). "Recent changes in extreme floods across multiple continents." *Environmental Research Letters* **12**(11): 114035.
- Blume, T., I. van Meerveld and M. Weiler (2018). "Incentives for field hydrology and data sharing: collaboration and compensation: reply to "A need for incentivizing field hydrology, especially in an era of open data"." *Hydrological Sciences Journal* **63**(8): 1266-1268.
- 680 [Blumer, A., A. Ehrenfeucht, D. Haussler and M. K. Warmuth \(1987\). "Occam's razor." \*Information processing letters\* \*\*24\*\*\(6\): 377-380.](#)
- Brodeur, Z. P., J. D. Herman and S. Steinschneider (2020). "Bootstrap Aggregation and Cross-Validation Methods to Reduce Overfitting in Reservoir Control Policy Search." *Water Resources Research* **56**(8): e2020WR027184.
- Buermann, W., J. Dong, X. Zeng, R. B. Myneni and R. E. Dickinson (2001). "Evaluation of the utility of satellite-based vegetation leaf area index data for climate simulations." *Journal of Climate* **14**(17): 3536-3550.
- 685 Center, G. G. R. D. (2005). Global Runoff Database, Koblenz.
- Ceola, S., B. Arheimer, E. Baratti, G. Blöschl, R. Capell, A. Castellarin, J. Freer, D. Han, M. Hrachowitz, Y. J. H. Hundecha and E. S. Sciences (2015). "Virtual laboratories: new opportunities for collaborative water science." *Hydrological Sciences Journal* **60**(4): 2101-2117.
- 690 Chagas, V. B., P. L. Chaffe, N. Addor, F. M. Fan, A. S. Fleischmann, R. C. Paiva and V. A. Siqueira (2020). "CAMELS-BR: hydrometeorological time series and landscape attributes for 897 catchments in Brazil." *Earth System Science Data* **12**(3): 2075-2096.
- China, M. o. G. a. M. R. o. t. P. s. R. o. (1991). "Geological map of Nei Mongol Autonomous Region, People's Republic of China, scale 1:1,500,000."
- Coron, L., V. Andreassian, C. Perrin, J. Lerat, J. Vaze, M. Bourqui and F. Hendrickx (2012). "Crash testing hydrological models in contrasted climate conditions: An experiment on 216 Australian catchments." *Water Resources Research* **48**(5).
- 695 Coxon, G., N. Addor, J. P. Bloomfield, J. Freer, M. Fry, J. Hannaford, N. J. Howden, R. Lane, M. Lewis and E. L. Robinson (2020). "CAMELS-GB: hydrometeorological time series and landscape attributes for 671 catchments in Great Britain." *Earth System Science Data* **12**(4): 2459-2483.
- Dai, Y., Q. Xin, N. Wei, Y. Zhang, W. Shangguan, H. Yuan, S. Zhang, S. Liu and X. Lu (2019). "A global high-resolution data set of soil hydraulic and thermal properties for land surface modeling." *Journal of Advances in Modeling Earth Systems* **11**(9): 2996-3023.
- 700 de Araújo, J. C. and J. I. González Piedra (2009). "Comparative hydrology: analysis of a semiarid and a humid tropical watershed." *Hydrological Processes: An International Journal* **23**(8): 1169-1178.
- Desborough, C. E. (1997). "The impact of root weighting on the response of transpiration to moisture stress in land surface schemes." *Monthly Weather Review* **125**(8): 1920-1930.
- 705 Didan, K. (2015). MOD13A3 MODIS/Terra vegetation indices monthly L3 global 1km SIN grid V006 (Data set) NASA EOSDIS Land Process, DAAC.
- Feng, D., K. Fang and C. Shen (2020). "Enhancing streamflow forecast and extracting insights using long-short term memory networks with data integration at continental scales." *Water Resources Research* **56**(9): e2019WR026793.



- Friedl, M. A., D. Sulla-Menashe, B. Tan, A. Schneider, N. Ramankutty, A. Sibley and X. Huang (2010). "MODIS Collection 5 global land cover: Algorithm refinements and characterization of new datasets." *Remote sensing of Environment* **114**(1): 168-182.
- 710 Gleeson, T., N. Moosdorf, J. Hartmann and L. Van Beek (2014). "A glimpse beneath ~~earth's~~ surface: GLobal HYdrogeology MaPS (GLHYMPS) of permeability and porosity." *Geophysical Research Letters* **41**(11): 3891-3898.
- Gleeson, T., L. Smith, N. Moosdorf, J. Hartmann, H. H. Dürr, A. H. Manning, L. P. van Beek and A. M. Jellinek (2011). "Mapping permeability over the surface of the Earth." *Geophysical Research Letters* **38**(2).
- 715 Gudmundsson, L., M. Leonard, H. X. Do, S. Westra and S. I. Seneviratne (2019). "Observed trends in global indicators of mean and extreme streamflow." *Geophysical Research Letters* **46**(2): 756-766.
- Hartmann, J. and N. Moosdorf (2012). "The new global lithological map database GLiM: A representation of rock properties at the Earth surface." *Geochemistry, Geophysics, Geosystems* **13**(12).
- 720 Hengl, T., J. Mendes de Jesus, G. B. Heuvelink, M. Ruiperez Gonzalez, M. Kilibarda, A. Blagotić, W. Shangguan, M. N. Wright, X. Geng and B. Bauer-Marschallinger (2017). "SoilGrids250m: Global gridded soil information based on machine learning." *PLoS one* **12**(2): e0169748.
- Horn, B. K. (1981). "Hill shading and the reflectance map." *Proceedings of the IEEE* **69**(1): 14-47.
- Huang, H., Y. Han, M. Cao, J. Song and H. Xiao (2016). "Spatial-temporal variation of aridity index of China during 1960–2013." *Advances in Meteorology* **2016**.
- 725 ~~Jenson, S. K., J. O. J. P. e. Domingue and r. sensing (1988). "Extracting topographic structure from digital elevation data for geographic information system analysis." *54*(11): 1593-1600.~~
- ~~Kendall, M. G. J. B. (1938). "A new measure of rank correlation." *30*(1/2): 81-93.~~
- Knoben, W. J., J. E. Freer and R. A. Woods (2019). "Inherent benchmark or not? Comparing Nash–Sutcliffe and Kling–Gupta efficiency scores." *Hydrology and Earth System Sciences* **23**(10): 4323-4331.
- 730 Knyazikhin, Y. (1999). "MODIS leaf area index (LAI) and fraction of photosynthetically active radiation absorbed by vegetation (FPAR) product (MOD 15) algorithm theoretical basis document." <http://eospsa.gsfc.nasa.gov/atbd/mod15atbd.html>.
- Kollat, J., P. Reed and T. Wagener (2012). "When are multiobjective calibration trade-offs in hydrologic models meaningful?" *Water Resources Research* **48**(3).
- Kollat, J., P. Reed and T. J. W. R. R. Wagener (2012). "When are multiobjective calibration trade-offs in hydrologic models meaningful?" **48**(3).
- 735 Kratzert, F., D. Klotz, ~~C. Brenner, K. Schulz and M. Herrnegger (2018). "Rainfall runoff modelling using long short term memory (LSTM) networks." *Hydrology and Earth System Sciences* **22**(11): 6005–6022.~~
- ~~Kratzert, F., D. Klotz, G. Shalev, G. Klambauer, S. Hochreiter and G. Nearing (2019). "Towards learning universal, regional, and local hydrological behaviors via machine learning applied to large-sample datasets." *Hydrology & Earth System Sciences* **23**(12).~~
- 740 Lane, R. A., G. Coxon, J. E. Freer, T. Wagener, P. J. Johnes, J. P. Bloomfield, S. Greene, C. J. Macleod and S. M. Reaney (2019). "Benchmarking the predictive capability of hydrological models for river flow and flood peak predictions across over 1000 catchments in Great Britain." *Hydrology and Earth System Sciences* **23**(10): 4011-4032.
- ~~LeCun, Y., Y. Bengio and G. J. Hinton (2015). "Deep learning." *521*(7553): 436–444.~~
- Legasa, M. and J. M. Gutiérrez (2020). "Multisite Weather Generators using Bayesian Networks: An illustrative case study for precipitation occurrence." *Water Resources Research* **56**(7): e2019WR026416.
- 745 Lehner, B. (2014). HydroBASINS: Global watershed boundaries and sub-basin delineations derived from HydroSHEDS data at 15 second resolution—Technical documentation version 1. c.
- Lehner, B., C. R. Liedmann, C. Revenga, C. Vörösmarty, B. Fekete, P. Crouzet, P. Döll, M. Endejan, K. Frenken and J. J. T. D. Magome, Version (2011). "Global reservoir and dam (grand) database." **1**: 1-14.
- 750 Linke, S., B. Lehner, C. O. Dallaire, J. Ariwi, G. Grill, M. Anand, P. Beames, V. Burchard-Levine, S. Maxwell and H. Moidu (2019). "Global hydro-environmental sub-basin and river reach characteristics at high spatial resolution." *Scientific data* **6**(1): 1-15.
- Liu, B., M. Xu, M. Henderson and W. Gong (2004). "A spatial analysis of pan evaporation trends in China, 1955–2000." *Journal of Geophysical Research: Atmospheres* **109**(D15).
- Liu, Q., Z. Yang and X. Xia (2010). "Trends for pan evaporation during 1959-2000 in China." *Procedia Environmental Sciences* **2**: 1934-1941.
- 755 Liu, Y., J. Zheng, Z. Hao and X. Zhang (2017). "Unprecedented warming revealed from multi-proxy reconstruction of temperature in southern China for the past 160 years." *Advances in Atmospheric Sciences* **34**(8): 977-982.
- ~~Maidment, D. R. (1996). *GIS and hydrologic modeling—an assessment of progress. Third International Conference on GIS and Environmental Modeling, Santa Fe, New Mexico.*~~
- ~~Maidment, D. R. and S. Morehouse (2002). *Arc Hydro: GIS for water resources, ESRI, Inc.*~~
- 760 Masutomi, Y., Y. Inui, K. Takahashi and Y. Matsuoka (2009). "Development of highly accurate global polygonal drainage basin data." *Hydrological Processes: An International Journal* **23**(4): 572-584.
- Mei, Y., V. Maggioni, P. Houser, Y. Xue and T. Rouf (2020). "A nonparametric statistical technique for spatial downscaling of precipitation over High Mountain Asia." *Water Resources Research* **56**(11): e2020WR027472.

- 765 Myneni, R., Y. Knyazikhin and T. Park (2015). "MOD15A2H MODIS/terra leaf area index/FPAR 8-day L4 global 500 m SIN grid V006." NASA EOSDIS Land Processes DAAC.
- Nevo, S., V. Anisimov, G. Elidan, R. El-Yaniv, P. Giencke, Y. Gigi, A. Hassidim, Z. Moshe, M. Schlesinger and G. Shalev (2019). "ML for flood forecasting at scale." arXiv preprint arXiv:1901.09583.
- Newman, A., M. Clark, K. Sampson, A. Wood, L. Hay, A. Bock, R. Viger, D. Blodgett, L. Brekke and J. Arnold (2015). "Development of a large-sample watershed-scale hydrometeorological data set for the contiguous USA: data set characteristics and assessment of regional variability in hydrologic model performance." Hydrology and Earth System Sciences **19**(1): 209-223.
- 770 Ni, H. and S. M. Benson (2020). "Using Unsupervised Machine Learning to Characterize Capillary Flow and Residual Trapping." Water Resources Research **56**(8): e2020WR027473.
- Oudin, L., V. Andréassian, J. Lerat and C. Michel (2008). "Has land cover a significant impact on mean annual streamflow? An international assessment using 1508 catchments." Journal of hydrology **357**(3-4): 303-316.
- 775 Running, S., Q. Mu and M. Zhao (2017). MOD16A2 MODIS/Terra Net Evapotranspiration 8-Day L4 Global 500 m SIN Grid V006. NASA EOSDIS Land Processes DAAC.
- Seybold, H., D. H. Rothman and J. W. Kirchner (2017). "~~Climate's~~Climate's watermark in the geometry of stream networks." Geophysical Research Letters **44**(5): 2272-2280.
- Shangguan, W., Y. Dai, Q. Duan, B. Liu and H. Yuan (2014). "A global soil data set for earth system modeling." Journal of Advances in Modeling Earth Systems **6**(1): 249-263.
- 780 Shangguan, W., Y. Dai, B. Liu, A. Zhu, Q. Duan, L. Wu, D. Ji, A. Ye, H. Yuan and Q. Zhang (2013). "A China data set of soil properties for land surface modeling." Journal of Advances in Modeling Earth Systems **5**(2): 212-224.
- Shen, C., E. Laloy, A. Elshorbagy, A. Albert, J. Bales, F.-J. Chang, S. Ganguly, K.-L. Hsu, D. Kifer and Z. Fang (2018). "HESS Opinions: Incubating deep-learning-powered hydrologic science advances as a community." Hydrology and Earth System Sciences (Online) **22**(11).
- 785 Silberstein, R. (2006). "Hydrological models are so good, do we still need data?" Environmental Modelling & Software **21**(9): 1340-1352.
- Singh, R., S. Archfield and T. Wagener (2014). "Identifying dominant controls on hydrologic parameter transfer from gauged to ungauged catchments—A comparative hydrology approach." Journal of Hydrology **517**: 985-996.
- Singh, R., K. van Werkhoven and T. Wagener (2014). "Hydrological impacts of climate change in gauged and ungauged watersheds of the Olifants basin: a trading-space-for-time approach." Hydrological Sciences Journal **59**(1): 29-55.
- 790 Subramanya, K. (2013). Engineering Hydrology, 4e, Tata McGraw-Hill Education.
- Sulla-Menashe, D. and M. A. Friedl (2018). "User guide to collection 6 MODIS land cover (MCD12Q1 and MCD12C1) product." USGS: Reston, VA, USA: 1-18.
- Survey, C. G. (2001). "1:2,500,000-scale digital geological map database of China."
- Tyralis, H., G. Papacharalampous and S. Tantanev (2019). "How to explain and predict the shape parameter of the generalized extreme value distribution of streamflow extremes using a big dataset." Journal of Hydrology **574**: 628-645.
- 795 van Werkhoven, K., T. Wagener, P. Reed and Y. J. W. R. R. Tang (2008). "Characterization of watershed model behavior across a hydroclimatic gradient." **44**(1).
- van Wijk, M. T. and M. Williams (2005). "Optical instruments for measuring leaf area index in low vegetation: application in arctic ecosystems." Ecological Applications **15**(4): 1462-1470.
- 800 ~~Vigliani, A., M. Borga, P. Balabanis and G. Blöschl (2010). "Barriers to the exchange of hydrometeorological data in Europe: Results from a survey and implications for data policy." Journal of Hydrology **394**(1-2): 63-77.~~
- Voepel, H., B. Ruddell, R. Schumer, P. A. Troch, P. D. Brooks, A. Neal, M. Durcik and M. Sivapalan (2011). "Quantifying the role of climate and landscape characteristics on hydrologic partitioning and vegetation response." Water Resources Research **47**(10).
- Wang, J., M. Chen, G. Lü, S. Yue, Y. Wen, Z. Lan and S. Zhang (2020). "A data sharing method in the open web environment: Data sharing in hydrology." Journal of Hydrology **587**: 124973.
- 805 Wickel, B., B. Lehner and N. Sindorf (2007). HydroSHEDS: A global comprehensive hydrographic dataset. AGU Fall Meeting Abstracts.
- Wongso, E., R. Nateghi, B. Zaitchik, S. Quiring and R. Kumar (2020). "A Data-Driven Framework to Characterize State-Level Water Use in the United States." Water Resources Research **56**(9): e2019WR024894.
- Woods, R. A. J. A. i. W. R. (2009). "Analytical model of seasonal climate impacts on snow hydrology: Continuous snowpacks." **32**(10): 1465-1481.
- 810 Xinjiang, B. o. G. a. M. R. o. (1992). "Geological map of Xinjiang Uygur, Autonomous Region, China, version 2, scale 1:1,500,000."
- Xu, Y., X. Gao, Y. Shen, C. Xu, Y. Shi and a. Giorgi (2009). "A daily temperature dataset over China and its application in validating a RCM simulation." Advances in Atmospheric sciences **26**(4): 763-772.
- Yamazaki, D., D. Ikeshima, J. Sosa, P. D. Bates, G. H. Allen and T. M. Pavelsky (2019). "MERIT Hydro: a high-resolution global hydrography map based on latest topography dataset." Water Resources Research **55**(6): 5053-5073.
- 815 Zeng, X. (2001). "Global vegetation root distribution for land modeling." Journal of Hydrometeorology **2**(5): 525-530.