

EStreams

An integrated dataset and catalogue of streamflow, hydro-climatic and landscape variables for Europe

do Nascimento, Thiago V.M.; Rudlang, Julia; Höge, Marvin; van der Ent, Ruud; Chappon, Máté; Seibert, Jan; Hrachowitz, Markus; Fenicia, Fabrizio

DOI

[10.1038/s41597-024-03706-1](https://doi.org/10.1038/s41597-024-03706-1)

Publication date

2024

Document Version

Final published version

Published in

Scientific Data

Citation (APA)

do Nascimento, T. V. M., Rudlang, J., Höge, M., van der Ent, R., Chappon, M., Seibert, J., Hrachowitz, M., & Fenicia, F. (2024). EStreams: An integrated dataset and catalogue of streamflow, hydro-climatic and landscape variables for Europe. *Scientific Data*, 11(1), Article 879. <https://doi.org/10.1038/s41597-024-03706-1>

Important note

To cite this publication, please use the final published version (if applicable).
Please check the document version above.

Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

Please contact us and provide details if you believe this document breaches copyrights.
We will remove access to the work immediately and investigate your claim.



OPEN

DATA DESCRIPTOR

EStreams: An integrated dataset and catalogue of streamflow, hydro-climatic and landscape variables for Europe

Thiago V. M. do Nascimento^{1,2,✉}, Julia Rudlang³, Marvin Höge¹, Ruud van der Ent³, Máté Chappon⁴, Jan Seibert², Markus Hrachowitz³ & Fabrizio Fenicia¹

Large-sample hydrology datasets have become increasingly available, contributing to significant scientific advances. However, in Europe, only a few such datasets have been published, capturing only a fraction of the wealth of information from national data providers in terms of available spatial density and temporal extent. We present “EStreams”, an extensive dataset of hydro-climatic variables and landscape descriptors and a catalogue of openly available stream records for 17,130 European catchments. Spanning up to 120 years, the dataset includes streamflow indices, catchment-aggregated hydro-climatic signatures and landscape attributes (topography, soils, geology, vegetation and landcover). The catalogue provides detailed descriptions that allow users to directly access streamflow data sources, overcoming challenges related to data redistribution policies, language barriers and varied data portal structures. EStreams also provides Python scripts for data retrieval, aggregation and processing, making it dynamic in contrast to static datasets. This approach enables users to update their data as new records become available. Our goal is to extend current large-sample datasets and further integrate hydro-climatic and landscape data across Europe.

Background & Summary

Large-sample datasets of hydrological variables across many catchments and long time periods are crucial for understanding and predicting hydrological variability in time and space^{1,2}. These datasets are increasingly in demand due to the rise of data-intensive machine learning models³.

Following the publication of the MOPEX dataset in the early 2000s, there has recently been a broad movement to making large-sample hydrology (LSH) datasets available. Many of those were developed inspired by the Catchment Attributes and Meteorology for Large-sample Studies (CAMELS) initiative that compiled and made available full datasets for the contiguous United States¹. Many countries and regions have embraced these or similar initiatives, including Australia⁴, Brazil⁵, Chile⁶, Great Britain², Switzerland⁷, Central-Europe⁸, North America⁹, China¹⁰, Central Asia¹¹ and Iceland¹².

At the global scale, there are already some collection efforts for hydro-meteorological data. The Global Streamflow Indices and Metadata Archive (GSIM)^{13,14} provides streamflow indices for 35,000+ locations around the globe, but no extensive set of catchment landscape and meteorological attributes. Recently another global streamflow indices time series initiative took place enlarging the analysis to 41,000+ river branches worldwide and using different streamflow signatures to enrich the flow regime analysis¹⁵. Considering streamflow records, the Global Runoff Data Centre (GRDC)¹⁶ provides data for 10,000+ stations, but similar to the previous datasets, no catchment attributes and meteorological forcing time series are available. In addition, the GRDC data is only updated episodically, while the others do, to our knowledge, not provide any updates. More recently the Caravan³ dataset compilation was published as a global initiative for standardizing already open-source

¹Eawag: Swiss Federal Institute of Aquatic Science and Technology, Dübendorf, Switzerland. ²Department of Geography, University of Zurich, Zurich, Switzerland. ³Department of Water Management, Faculty of Civil Engineering and Geosciences, Delft University of Technology, Delft, Netherlands. ⁴Széchenyi István University, Department of Transport Infrastructure and Water Resources Engineering, Győr, Hungary. ✉e-mail: thiago.nascimento@eawag.ch

published streamflow datasets of initially 6,830 catchments, where catchment attributes and meteorological forcing were derived from gridded global products.

While global datasets offer easy access, they come with limitations. Firstly, their spatial coverage remains restricted, offering only a fraction of data available from national providers worldwide. The Caravan dataset, for example, originally covered Europe for only Great Britain, Austria and the Danube catchment as far downstream as the city of Bratislava (Slovakia). By now, there are multiple extensions for Denmark, Israel, Switzerland, Spain, Iceland and, most recently, a GRDC extension¹⁷ adding another 25 countries globally. Yet, for eastern and southern Europe publicly available data is still difficult to access. Secondly, such datasets are also limited in their temporal extent. For example, the CAMELS-GB² covers the period from 1970 to 2015, while the LamaH-CE dataset⁸ spans from 1981 to 2017. Thirdly, existing large sample hydrology datasets, including the CAMELS databases, lack extensibility, making the accommodation of newly available data challenging.

Although most countries collect daily streamflow data at numerous river gauging stations, compiling a comprehensive hydrological dataset from this information presents significant challenges. Firstly, access to these data can be challenging. Some countries offer this data on the official websites of government agencies or associated data providers, while others provide it upon request. Official government websites are frequently available only in national languages, adding an extra layer of complexity. Gaining access can be intricate, involving navigation to a selection of stations and periods, which need to be downloaded individually. Secondly, substantial formatting and pre-processing are often necessary before the data can be effectively utilized. Finally, redistribution restrictions may hinder the republishing of country-specific data. These obstacles pose significant barriers to hydrological analyses of catchments in large-sample investigations, particularly given the short timeframes of typical research projects.

Here, we present “EStreams”, a platform consisting of two distinct products: (1) an extensive streamflow catalogue together with Python scripts for data direct access at the individual data providers and (2) a dataset of weekly, monthly, seasonal and annual indices, of streamflow, together with the associated catchment-averaged hydro-climatic signatures, meteorological time series and landscape descriptors for 17,130 catchments across 41 countries over pan-European territory. Currently, the dataset covers the period of 1900–2022.

While the focus of EStreams is on streamflow, the EStreams dataset also contains catchment aggregated meteorological forcing and landscape descriptors, typically necessary for hydrological analyses. These indices and descriptors were derived from various open source datasets and include climate¹⁸, geology^{19,20}, hydrology and topography^{21–24}, land use and land cover^{25–27}, soil types^{28–30} and vegetation characteristics^{31,32}. Similarly to streamflow, national providers often have more accurate information for such auxiliary data, but seldom they are easily accessible.

Unlike existing global datasets, which are relatively “static” as not easily updatable with new stations or recent time periods, EStreams is designed as “dynamic” by linking users to the original data providers. While “static” datasets may offer more accurate quality checks and are well-suited for applications such as benchmarking methods and models, many practical applications benefit from using the most up-to-date and dense data. This is particularly true for tasks like accurate streamflow predictions using data-intensive machine learning models.

Hence, our main contributions with this work are:

- i. Introducing the currently most extensive and extensible integrated collection of weekly, monthly, seasonal and annual indices of streamflow for Europe, along with catchment-aggregated meteorological and landscape variables (dataset).
- ii. Providing detailed metadata for streamflow gauges, including catchment boundaries, and a catalogue of the corresponding data providers.
- iii. Allowing reproducibility and extension by making available all codes used to retrieve the source data and aggregate them by catchment in an easy-to-use workflow, allowing users to directly and readily access the desired data from data providers.

The methodology employed to process the source data and obtain the current dataset and catalogue is illustrated in Fig. 1. This figure highlights the primary data sources, the general procedure, and the final outputs of EStreams. A detailed description of each step is provided in the Methods sections.

Methods

Streamflow data. *Available stations.* Daily streamflow data from 17,130 European river catchments with varying sizes and characteristics were aggregated from 41 countries and more than 50 different data providers. In some countries, such as Italy and Germany, multiple data providers contributed to the dataset. Figure 2a shows the distribution of the gauges with their respective catchment boundaries in the background. As can be seen in the figure, there is a significant variability in terms of station density, which is the highest in central Europe and the lowest in the South and the East. The time series records span the period 1900–2022, with varying length for each catchment, as shown in Fig. 2b. Central Europe features the longest time series, with many stations with records extending over 80 years. Figure 2c shows the evolution of the number of stations with measurements at a given time accounting for the discontinuity of stations over time. The plot shows an increasing trend in the number of gauging stations with concurrent records.

The streamflow records were selected based on the following criteria: (i) they were available from official authorities in their respective country or from a recent open-access dataset, and (ii) they were open-source and easily accessible either via the internet or by e-mail request. The latter point emphasizes that no dataset requiring purchase for non-commercial access were included. It is important to note that freely available data do not necessarily come with a free redistribution license. Therefore, we cannot and do not make raw daily streamflow data directly available. Should the source data be necessary, we provide the EStreams catalogue of data sources

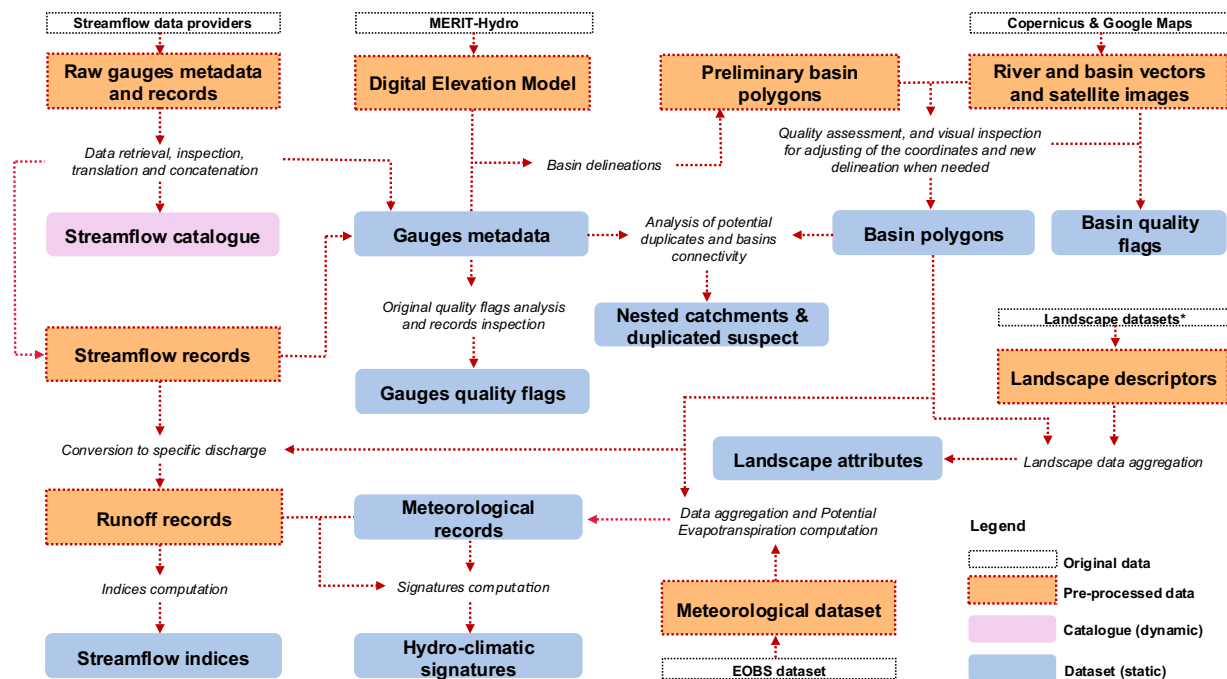


Fig. 1 Framework of the methodology adopted in EStreams for deriving the Streamflow Catalogue, and the Dataset. The boxes with dashed lines represent the original, and the intermediate (pre-processed) data used in EStreams. The outputs are shown in pink (catalogue) and blue (dataset). *The landscape datasets encompass topography, soils, geology, hydrology, vegetation and land cover.

to allow users easy and direct data access from the original repositories, including codes and instructions for data download and formatting. Compared to static databases of pre-compiled datasets currently available, our approach has two main advantages:

- i. Users can tailor the download to determine the desired spatial and temporal coverage, also making use of the provided descriptive statistics of the source data, such as regime characteristics or catchment properties.
- ii. Users can access the most up-to-date information directly from the data sources.

Table 1 provides an overview of the contributing countries, the number of streamflow gauges, and the data providers. France has the highest number of gauges (4,968), followed by Germany (2,093) and Spain (1,440). In contrast, Bulgaria (8 gauges) Moldova (2) and North Macedonia (1) have the lowest numbers of gauges.

Streamflow gauges labelling. After the collection of the streamflow data and gauge information from each provider, the individual datasets were collated into a single dataset. In this process, each gauge was labelled with a unique 8-digit code. Consequently, each catchment was renamed according to its respective streamflow gauge. The 8-digit codes were generated using the following logic: the first two digits represent the country/region, the next two digits represent specifications about the data provider within regions that had more than one official provider, and the last four digits refer to the gauge counter for each country/region. For example, the gauge GB000045 represents Great Britain (GB), with only one provider (00), and the gauge number 0045. Similarly, ITIS0001 represents Italy (IT), with ISPRA (IS) as the data provider, and gauge number 0001. The gauges with records obtained from GRDC have the second two digits as “GR” (e.g., LVGR0001) to facilitate identification. This standardization ensures that all gauges are consistently labelled, providing users with a clear indication of the source and the number of records.

Identification of duplicate gauges. When compiling large streamflow datasets, there is a possibility of having duplicate records within the dataset that need to be identified and removed. This issue can arise when combining information from multiple sources and even within datasets obtained from a single data provider. To identify suspected duplicate records, we used a similar approach as used by the GSIM¹³, where for gauges originating from distinct data providers, we identified potential duplicate gauges by examining similarities in gauge and river names. We employed the Jaro-Winkler distance metric to quantify alphanumeric similarity, as discussed by Christen, 2012³³ with a threshold set at 0.70. We additionally considered spatial proximity, constraining pairs of stations within 1 km of each other. For gauges originating from the same data provider, we selected stations within a spatial proximity of 50 m and a delineated area difference below 1%. Gauges meeting these criteria were flagged as potential duplicates. The list of potential duplicates for each gauge is contained in the attribute *duplicated_suspect* within the gauges’ layer in the final EStreams dataset. Notably, all potential duplicates are

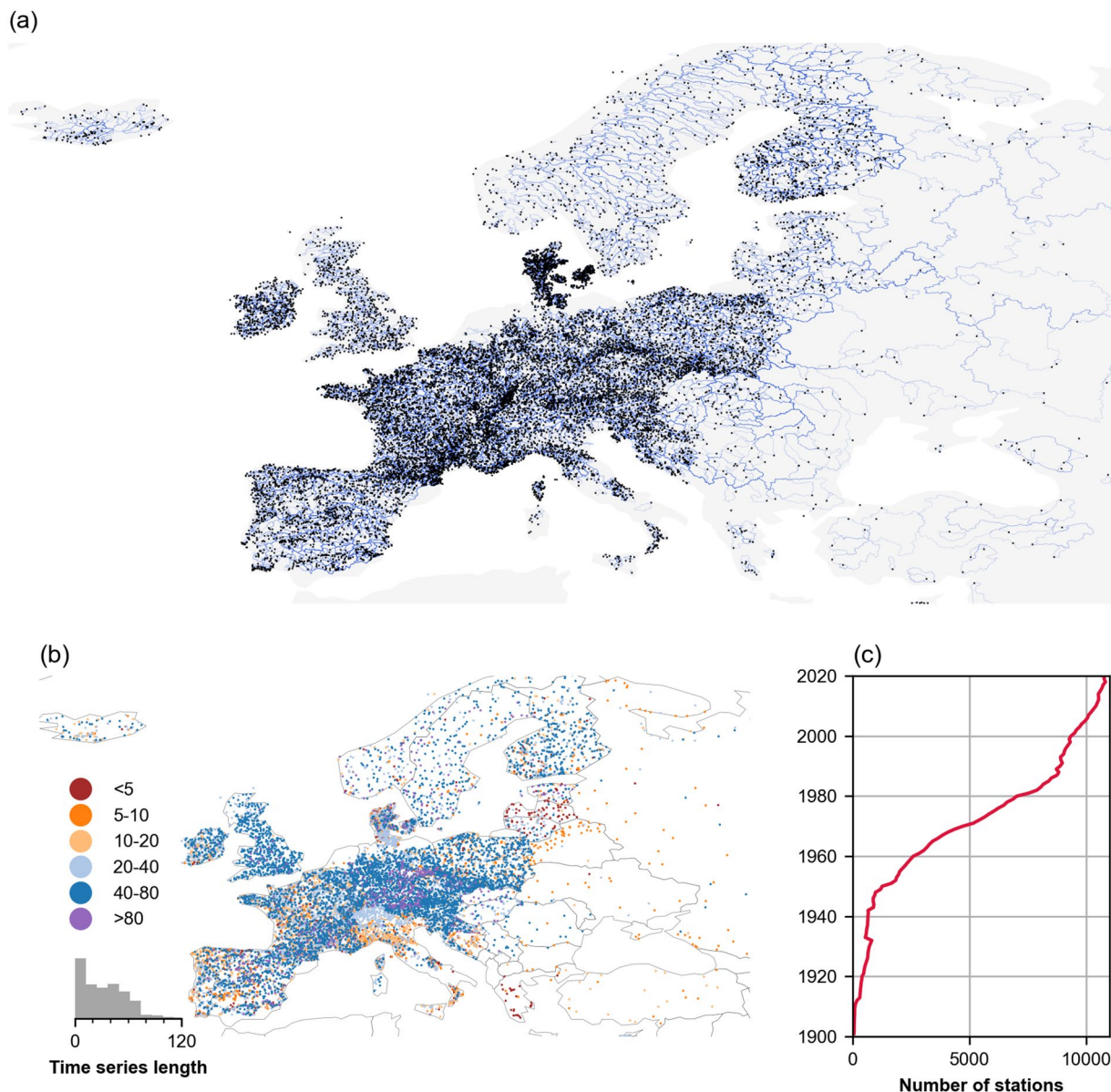


Fig. 2 (a) Spatial distribution of the 17,130 streamflow gauges currently included in EStreams (in black dots) with their catchment boundaries in background (in blue) over Europe. (b) Spatial distribution of the streamflow with the colors representing the time series length in years. (c) Temporal evolution of station coverage. The plot shows the number of active stations in a given year, Although the curve accounts for dismissed stations, it still shows an increasing trend. Basemap from GeoPandas¹⁰⁴.

preserved in EStreams, giving users the flexibility to choose their preferred station and data provider when duplicates are found. This approach ensures that users can tailor their dataset according to their specific needs and preferences.

Quality flags of records. Quality control of streamflow data is essential before undertaking any hydrological study. While some data providers include quality flags with each published record, this practice is not consistently available. Automatic checks are available but may be subjective, and their effectiveness has not yet been fully investigated^{34,35}. For example, Do, 2018¹³ employed an automatic detection criterion to identify and filter potentially suspect records based on negative values, consecutive repetitions, and outliers. However, these filtering criteria are not always reliable, as pointed out by Chen, 2023¹⁵.

In this work, following the approach utilized by Chen, 2023¹⁵, we adopt a two stages approach for quality checking the data, the first oriented at individual data points, and the second assessing the entire record. The first stage is primarily based on the quality flags from the original providers, when available, which for consistency are reclassified into four categories: “missing”, “no-flags”, “suspect” and “reliable”. First, all negative values were replaced with “not a number” (NaN) and flagged as “missing”. Then, values with a quality flag given by the data

| Country/region | Code | Stations | References |
|----------------|------|----------|--|
| Austria | AT | 582 | BML ⁴⁹ |
| Bosnia and H. | BA | 91 | GDRC ¹⁶ ; FHMZBIH ⁵⁰ |
| Belgium | BE | 230 | VW ⁵¹ ; SPW ⁵² |
| Bulgaria | BG | 8 | GRDC ¹⁶ |
| Belarus | BY | 51 | GRDC ¹⁶ |
| Switzerland | CH | 298 | BAFU ^{7,53} |
| Cyprus | CY | 14 | GRDC ¹⁶ |
| Czechia | CZ | 566 | CHMI ⁵⁴ |
| Germany | DE | 2,093 | LHW ⁵⁵ ; ASOEG ⁵⁶ ; Umweltportal ⁵⁷ ; ELWAS-WEB ⁵⁸ ; NLWKN ⁵⁹ ; HLNUG ⁶⁰ ; GKD ⁶¹ ; LUBW ⁶² ; WB ⁶³ ; LBAW ⁶⁴ ; MKUEM ⁶⁵ ; LUBN ⁶⁶ ; BFG ⁶⁷ |
| Denmark | DK | 1,000 | ODA ⁶⁸ |
| Estonia | EE | 67 | GRDC ¹⁶ |
| Spain | ES | 1,440 | CEDEX ⁶⁹ |
| Finland | FI | 669 | FEI ⁷⁰ |
| France | FR | 4,968 | BanqueHydro ⁷¹ |
| Great Britain | GB | 671 | NRFA ⁷² |
| Greece | GR | 31 | GRDC ¹⁶ ; OHIN ⁷³ ; HCRM ⁷⁴ |
| Croatia | HR | 317 | DHZ ⁷⁵ |
| Hungary | HU | 98 | GRDC ¹⁶ ; OVF ⁷⁶ |
| Ireland | IE | 464 | EPA ⁷⁷ ; OPW ⁷⁸ |
| Iceland | IS | 111 | LamaH-Ice ¹² |
| Italy | IT | 767 | GRDC ¹⁶ ; ISPRA ⁷⁹ ; APC Abruzzo ⁸⁰ ; CFRA Valle d'Aosta ⁸¹ ; ARPAE Emilia-Romagna ⁸² ; ARPA: Umbria ⁸³ ; Sardegna ⁸⁴ ; Lombardia ^{85,86} ; Toscana ⁸⁷ ; Piemonte ⁸⁸ ; ARPAL Liguria ⁸⁹ ; ARPAV Veneto ⁹⁰ ; SPRUD Trentino ⁹¹ |
| Lithuania | LT | 76 | GRDC ¹⁶ |
| Luxembourg | LU | 19 | NGGL ⁹² |
| Latvia | LV | 61 | GRDC ¹⁶ |
| Moldova | MD | 2 | GRDC ¹⁶ |
| Macedonia | MK | 1 | GRDC ¹⁶ |
| N. Ireland | NI | 51 | NRFI ⁷² |
| Netherlands | NL | 17 | RWS ⁹³ |
| Norway | NO | 189 | NVE ⁹⁴ |
| Poland | PL | 1,287 | IMGW-PIB ⁹⁵ |
| Portugal | PT | 280 | SNIRH ⁹⁶ |
| Romania | RO | 18 | GRDC ¹⁶ |
| Serbia | RS | 18 | GRDC ¹⁶ |
| Russia | RU | 98 | GRDC ¹⁶ |
| Sweden | SE | 290 | SMHI ⁹⁷ |
| Slovenia | SI | 117 | ARSO ⁹⁸ |
| Slovakia | SK | 21 | GRDC ¹⁶ |
| Turkey | TR | 28 | GRDC ¹⁶ |
| Ukraine | UA | 21 | GRDC ¹⁶ |

Table 1. Overview of streamflow time series data available per country/region, with information about number of stations and data providers.

providers had their original labels reclassified as either “reliable”, “suspect” or “missing”. Finally, all data without a quality flag from the original providers were classified as “no-flag”. A complete overview of the mapping between the original flags and our four flags system is available in Supplementary Table 1.

In the second stage, we assessed the overall reliability of each entire time series based on the fraction of problematic data points as determined in the previous stage. This classification considered five criteria outlined in Table 2.

A total of 7,430 stations had quality flags from their providers (about 43% of the total). Figure 3a shows that approximately 134 million data points (63.4% of the total) were classified as “no-flag”, 56 million data points (26.7%) as “reliable”, 3.9 million data points (1.9%) as “suspect”, and 16.8 million data points (8%) as “missing”. Regarding the gauge’s quality classification, Fig. 3b shows that most stations were categorized as either Class A or B (9,652), followed by Class E (3,317), Class C (2,827) and Class D (1,334). This classification allows users to filter the data depending on their needs. It is noteworthy that many national providers may offer only high-quality data for download. Therefore, even without explicit quality flags, the data can often be assumed to come from reliable stations. The quality flag for each gauge’s records is stored as the attribute *gauge_flag* within the gauges’ layer in the final EStreams dataset.

| Quality flag (gauge) | Criterion |
|----------------------|--|
| A | More than 95% of the gauge records flags are “reliable” |
| B | More than 95% of the gauge records flags are “reliable” or “no-flag” |
| C | Less than 10% of the gauge records flags are “missing” |
| D | Less than 20% of the gauge records flags are “missing” |
| E | More than 20% of the gauge records flags are “missing” |

Table 2. Criteria used for the quality assessment of the streamflow gauges as in Chen, 2023¹⁵. When one station met multiple criteria simultaneously, the highest-level flag was applied.

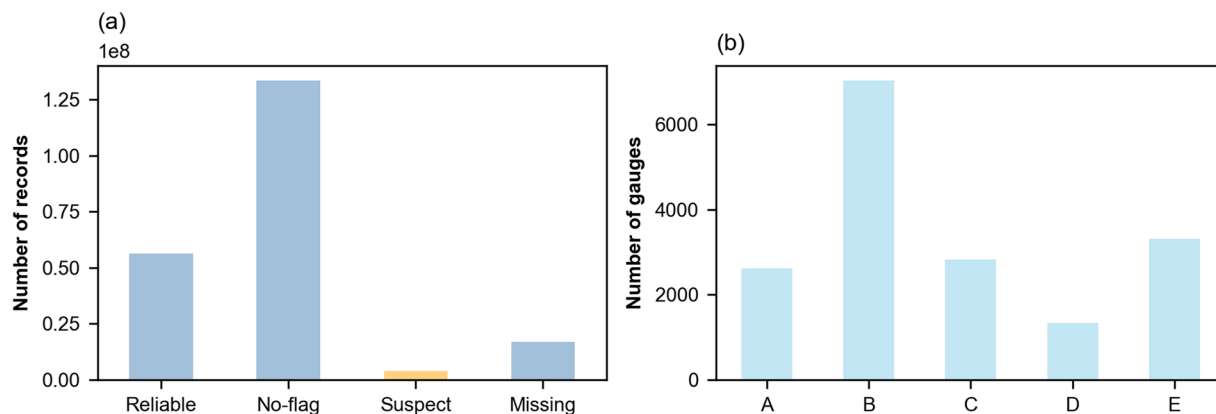


Fig. 3 (a) Histogram of the streamflow data points according to their four data quality flags and (b) Histogram of the number of gauges according to their integrated data quality flag.

Basin delineation. Since catchment boundaries shapefiles were rarely available from national providers, this work adopted a semi-automatic delineation of catchment boundaries corresponding to streamflow gauges using Python scripts and QGIS software. We used the “delineator” python package³⁶, which determines catchment boundaries using hybrid vector and raster-based methods. This package requires as input the latitude and longitude coordinates of the streamflow gauges and uses the MERIT-Hydro Digital Elevation Model (DEM)²¹. MERIT-Hydro is a digital elevation model developed to remove multiple error components from the existing spaceborne DEMs (SRTM3 v2.1 and AW3D-30m v1).

To appraise the accuracy of the delineated area, catchments were split into two categories: (i) catchments with a reported area from the data providers and (ii) catchments without this information. For gauges with available official catchment areas, the reported area was compared to the derived area, and the following workflow was adopted:

- i. First, we computed the “relative area difference” A_{rel} as defined in Eq. 1. If $|A_{rel}|$ was below 10%, regardless of catchment size, the delineation was accepted, and the catchment was labelled with a quality flag of “0”.
- ii. Otherwise, the catchment delineation was visually inspected, potentially corrected as described below, and assigned a specific quality flag as detailed in Table 3, which provides an overview of the flags used and number of gauges corresponding to each flag.

$$A_{rel} = 100 \times \frac{A_{EStream} - A_{official}}{A_{official}} \quad (1)$$

where $A_{EStream}$ is the calculated area in EStreams and $A_{official}$ is the reported official area.

The visual inspection was made using the river networks from both the MERIT-Hydro and EU-Hydro datasets³⁷, Google Maps satellite imagery, and nearby catchments delineated and labelled with a quality flag of “0”. These three data sets were used as they represent independent sources and offer a good trade-off for evaluating the catchment delineation usability.

During the visual inspection, it was observed that some boundary discrepancies could be corrected with an adjustment in the streamflow gauge location. We assumed that uncertainties in the georeferenced system or the presence of close-by river branches could cause these discrepancies. For those catchments, the gauge location was moved (snapped) to the closest point within the MERIT-Hydro River network based on the gauge’s river and location names.

Catchments with $|A_{rel}|$ below 10% after the snap were labelled with a quality flag “1” indicating accepted delineation after the snap. The remaining catchments were classified with the criteria detailed in Table 3.

| Basin area quality flag | Number of gauges | Description |
|-------------------------|------------------|--|
| 0 | 12,801 | $ A_{rel} $ below 10%. |
| 1 | 164 | $ A_{rel} $ below 10% after moving the gauge location. |
| 2 | 1,037 | $ A_{rel} $ above 10% or no reported area available, but delineation visually compared to other delineations from down and upstream gauges labelled “0”; Google Maps satellite imagery and to the EU-Copernicus River network. |
| 3 | 369 | $ A_{rel} $ above 10% or no reported area available, but delineation visually compared to Google Maps satellite imagery and to the EU-Copernicus River network. |
| 4 | 343 | $ A_{rel} $ above 30% or no reported area available, but delineation compared to EU-Copernicus River network. |
| 5 | 68 | $ A_{rel} $ above 10% or no reported area available, and delineation manually adjusted using EU-Copernicus in addition to MERIT-Hydro. |
| 6 | 11 | Similar to “5”, but still with $ A_{rel} $ above 30% or no reported area available. |
| 888 | 64 | $ A_{rel} $ above 10% or no reported area available, but location in areas under high human influence, such as canalization and water exports and in karstic regions. |
| 999 | 2,273 | $ A_{rel} $ above 10% or no reported area available, and delineation eventually not accepted after visual inspection. |

Table 3. Description of the catchment area quality flags adopted for the current catchment delineations and overview of the number of catchments per group.

It is important to note that for some situations where human-influence such as canalization, water exports and specific lithologies like karstic systems, the actual catchment boundary delineation remains challenging. Hence, for catchments where $|A_{rel}|$ was above 10% and the visual inspection indicated such situations, we assigned a quality flag of “888”.

Finally, catchments where $|A_{rel}|$ was above 10%, and were not visually adjusted or accepted, were assigned to a quality flag “999”.

Out of a total of 17,130 stations, 15,775 (92%) had a reported catchment area from the data providers. Figure 4a shows the distribution of these streamflow gauges divided into two classes: gauges with $|A_{rel}|$ above 50% (in red), and those with $|A_{rel}|$ below 50% (in blue). Generally, gauges with high area discrepancies are located in regions of low relief, partly canalized landscapes and with high presence of lakes such as in Denmark, Sweden and Croatia.

Figure 4b shows the exceedance percentage of $|A_{rel}|$ of these 15,775 catchments with a reported area. As indicated with the dashed orange line, the catchments with $|A_{rel}|$ above 50% was 8% (1,205 catchments). This analysis also shows that less than 17% of the catchments (2,712) had $|A_{rel}|$ above 10%.

Figure 4c focuses on catchments with $|A_{rel}|$ above 50% (1,205 catchments) and shows how the fraction of these catchment varies with catchment area. Notably, 17% of catchments under 100 km² exhibited $|A_{rel}|$ above 50%, while in all other ranges shown in the bar plot, the occurrence was below 5%. This analysis suggests that catchments with significant area differences tend to be relatively small.

Finally, for the 1,355 gauges (8% of the data) without catchment area information, the delineation was visually inspected, and a label was assigned to indicate the accuracy of the delineation based on the criteria shown in Table 3. Note that as it is not possible to calculate $|A_{rel}|$ for these catchments, the quality flags of “0” or “1” were never assigned to such basins. The visual inspection was again made using the river name, the river network provided by MERIT-Hydro and the EU-Hydro, Google Maps satellite imagery and nearby catchments delineated and labelled with a quality flag of “0”.

Hence, in the gauges’ layer stored in the final EStreams dataset, besides the original *lat* and *lon* coordinates, we included the *lat_snap* and *lon_snap* coordinates after the potential snap. The gauges layer also received an attribute called *area_estreams*, which express the $A_{EStream}$. Additionally, we included the A_{rel} as the attribute *area_rel*, and the qualitative flag as the attribute *area_flag*.

Catchment aggregated data. The EStreams dataset includes streamflow, meteorological, and landscape variables. For streamflow, we distinguish between dynamic streamflow indices and hydro-climatic signatures, which are further detailed in their respective sections. Meteorological variables are discussed in the “Meteorological records” section. Finally, landscape attributes were categorized into six groups (Topography, Soils, Geology, Hydrology, Vegetation, and Land Cover) and are described in the “Landscape attributes” section. All catchment aggregations were derived using the catchment boundaries and areas calculated by EStreams. For example, all streamflow indices and signatures were computed using the specific discharge (in mm/day) derived with the $A_{EStreams}$ areas.

Streamflow indices. In EStreams, streamflow data is presented in terms of “indices”, hence statistics of the daily data such as mean streamflow, maximum, minimum, percentiles and coefficient of variation, which are provided at annual, seasonal, monthly and weekly resolutions. The use of these indices is consistent with earlier works, such as the GSIM dataset^{13,14} and the CCI/WCRP/JCOMM Expert Team on Climate Change Detection and Indices (ETCCDI) (<https://www.wcrp-climate.org/data-etccdi>).

The use of indices instead of the daily data allows to make relevant climate information publicly available in cases where access to raw daily values is restricted. The selected indices, as discussed in the GSIM dataset^{13,14}, are

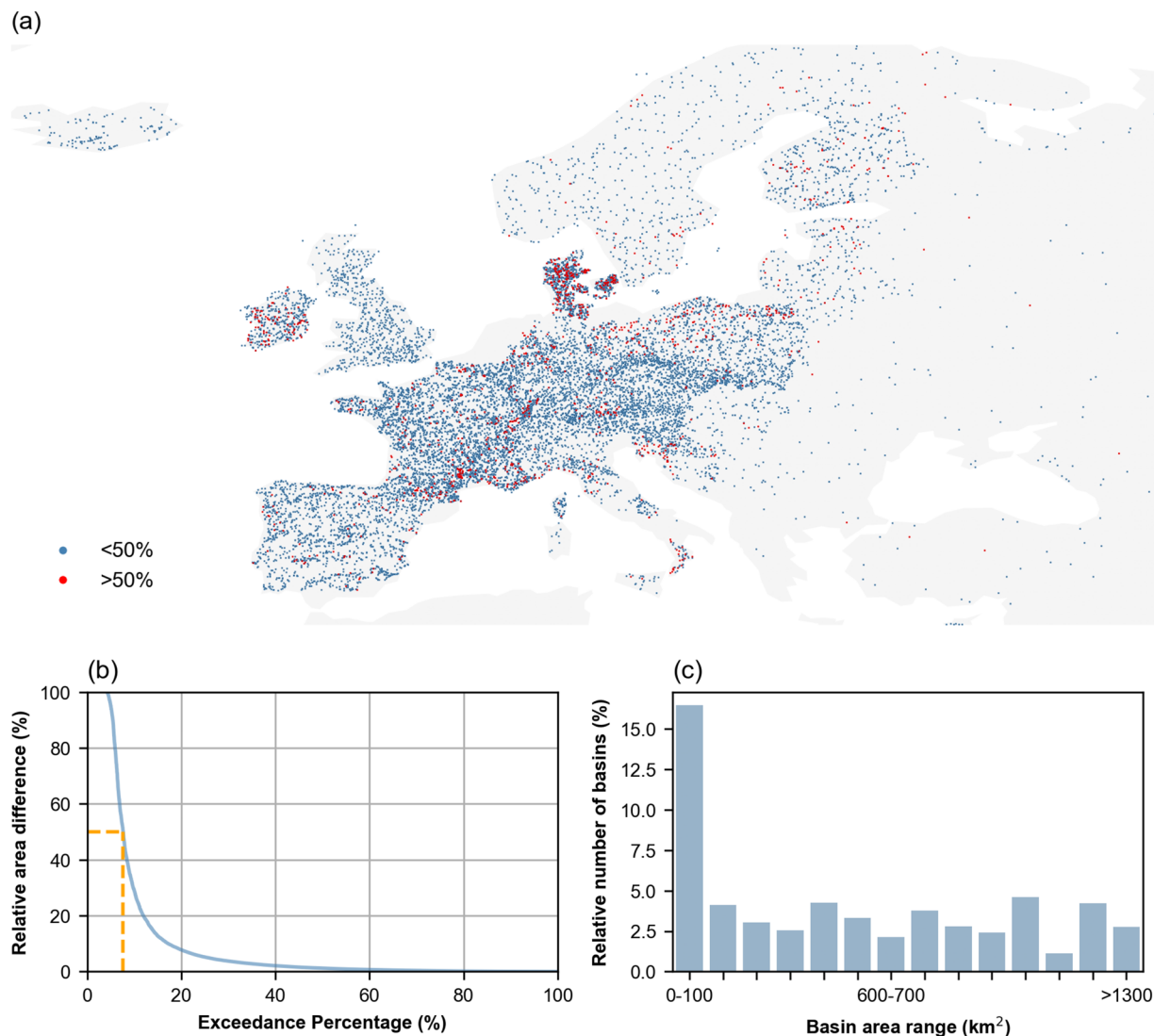


Fig. 4 (a) Relative absolute area difference $|A_{rel}|$ above 50% (in red) and below 50% (in blue). (b) Exceedance percentage of the $|A_{rel}|$; the orange line marks the exceedance percentage corresponding to a $|A_{rel}|$ of 50%. (c) Bar plots showing the relative number of basins with areas above 50% for different basin area ranges (e.g., 0–100 km², 100–200 km², and >1,300 km²) relative to the total number of basins in each range. Basemap from GeoPandas¹⁰⁴.

of high relevance and have been widely used in many hydrological studies, as they can facilitate the analysis of trends and changes in the regional water balance and the seasonal cycle.

The streamflow indices contained in EStreams are presented in Table 4, alongside with their units and temporal resolution. All the indices were computed for time-steps where at least 95% of the data was available, e.g., at annual time-step, the indices were computed for years where at least 347 days of data were available.

Hydro-climatic signatures. In addition to the streamflow indices, we computed the same set of meteorological and hydrological signatures provided in the original CAMELS dataset¹. Unlike streamflow indices, these signatures were calculated for the entire time period between 1950–2022 where data are available. Here we refer to these indices and signatures as hydro-climatic signatures (e.g., streamflow & precipitation mean, seasonality & aridity index, and runoff coefficient). For meteorology, we used precipitation and temperature derived from the Ensembles Observation (E-OBS) product¹⁸. This work used the “hydroanalysis” python package³⁸ for the computation of these signatures.

The full list of signatures used is available in Table 5. We considered only catchments with more than one year of continuous measurements within the period of 1950–2022. Additionally, we also provide the number of years used for the signature’s computation (*num_years*), the start (*start_date*) and the end (*end_date*) of the observations between 1950–2022 to give a further overview of the period the signature refers to, considering separately the hydrological (*hydro*) and the climatic (*climatic*) signatures.

| Variable | Description | Units | Resolution |
|--|---|----------------------|---------------|
| mean | Mean daily streamflow. | mm day ⁻¹ | W, M, S and Y |
| std | Standard deviation of the daily streamflow. | mm day ⁻¹ | W, M, S and Y |
| cv | Coefficient of the variation of the daily streamflow. | — | W, M, S and Y |
| min | Minimum daily streamflow. | mm day ⁻¹ | W, M, S and Y |
| max | Maximum daily streamflow. | mm day ⁻¹ | W, M, S and Y |
| min7 | Minimum 7-day streamflow. | mm day ⁻¹ | M, S and Y |
| max7 | Maximum 7-day streamflow. | mm day ⁻¹ | M, S and Y |
| p_{10, 20, 30, 40, 50, 60, 70, 80, 90} | Percentile values of the daily streamflow. | mm day ⁻¹ | S and Y |
| iqr | Interquartile range of the daily streamflow (P75 minus P25) | mm day ⁻¹ | W, M, S and Y |
| ct | Centre timing, which corresponds to the day of the year (doy) at which 50% of the annual flow is reached. | day | Y |
| doymin | The day of the year (doy) at which the minimum streamflow occurred. | day | Y |
| doymax | The day of the year (doy) at which the minimum streamflow occurred. | day | Y |
| doymin7 | The day of the year (doy) at which the minimum 7-day streamflow occurred. | day | Y |
| doymax7 | The day of the year (doy) at which the maximum 7-day streamflow occurred. | day | Y |
| gini | Gini coefficient | — | Y |

Table 4. Set of dynamic streamflow time series indices computed and made available at the present dataset.

Meteorological records. EStreams used E-OBS¹⁸ for meteorological forcing data records, which has been widely used in hydrological studies over Europe^{39–42}. E-OBS provides a pan-European observational dataset of surface climate variables that is derived by statistical interpolation of *in-situ* measurements, collected from national data providers. It is an open-access database with daily records ranging from 1950–present. We used the ensemble mean dataset at a resolution of 0.25 degrees. Additionally, we used the temperature records from E-OBS to derive potential evapotranspiration (PET) using the Hargreaves formulation⁴³ and the “pyet” python package⁴⁴ for computation. Each catchment has 9 daily meteorological time series associated with it, which are illustrated in Table 6. The accuracy of E-OBS may be dependent on station density⁴², which varies across Europe. In order to account for this potential source of uncertainty, EStreams also includes information on the number of weather stations and density aggregated to a buffer of 10 km within each catchment boundary.

Landscape attributes. A full overview of the landscape attributes contained in EStreams is shown in Table 7 and Table 8, with a short description, their units, and data provider. Regarding spatial coverage, except for the landcover & land use and soil types that have pan-European coverage, all the remaining products are global. Table 7 covers solely the fully static attributes, which are considered time invariant, such as elevation, soil types, main geology and mean vegetation indices. Conversely, Table 8 encompasses a group of attributes that are considered time variable, such as normalized difference vegetation index (NDVI), leaf-area index (LAI), irrigation and snow cover. These attributes are reported in time series at either monthly, yearly or in a specific number of years (e.g., irrigation and landcover) resolution.

Topographical attributes were based on MERIT-Hydro²¹. Geology made use of the widely used Global Lithological Map Database (GLiM)¹⁹ and a gridded product for the estimation of the depth to bedrock²⁰, which have been both used in several applications databases^{1,8,23}. For the number of dams and of total upstream reservoir volume we used the Georeferenced global dams and reservoirs dataset²². A similar aggregation was performed for lakes using the HydroLakes dataset⁴⁵. Vegetation indices and snow cover percentage made use of three MODIS products^{27,31,32} and were aggregated considering both temporal and static attributes. For irrigation, we decided to use the global dataset of the extent of irrigated land²⁶, which ranges from 1900 to 2005, and has been already used in other studies^{13,14,23}. The soil attributes were based on the European Soil Database Derived data (ESDD)^{28,29,30} and the land cover on the CORINE land cover dataset²⁵. Both are widely used products which have been used in previous LSH datasets covering Europe^{7,8}.

Data Records

The current version of the EStreams dataset and catalogue (v1.0) is stored at a Zenodo repository⁴⁶ at <https://doi.org/10.5281/zenodo.13154470>. The repository is organized into the following subfolders:

- **streamflow_gauges:** Contains two csv-files. One includes all the metadata associated with each of the 17,130 streamflow gauging stations such as location, river name, catchment area, and gauge elevation. The other file is the streamflow catalogue containing all the data provider information, further described in the following section.
- **shapefiles:** Contains two shapefiles. One shapefile includes the derived catchment boundaries associated with each streamflow gauge, and the other shapefile marks the location of the streamflow gauges. Both files are referenced in WGS 84.
- **streamflow_indices:** Contains one sub-folder per time resolution (weekly, monthly, seasonal and yearly) with a csv-file per computed index. The rows of each csv-file represent the time, and the columns represent the catchment.

| Signature | Unit | Description |
|------------------------------|-----------------------|--|
| q_mean | mm day ⁻¹ | Mean daily streamflow. |
| runoff_ratio | — | Ratio of mean daily streamflow to mean daily precipitation computed using Eq. (2) in Sawicz, 2011 ¹⁰⁰ . |
| q_elas_Sankarasubramanian | — | Streamflow precipitation elasticity. It represents the sensitivity of streamflow to changes in precipitation at the annual timescale computed using Eq. (7) in Sankarasubramanian, 2001 ⁹⁹ , the last element being P/Q not Q/P |
| slope_sawicz | — | Slope of the flow duration curve computed using Eq. (3) in Sawicz, 2011 ¹⁰⁰ . |
| baseflow_index | | Ratio of mean daily baseflow to mean daily streamflow. Hydrograph separation performed using the Ladson, 2013 ¹⁰¹ digital filter. |
| hfd_mean | day of year | Mean half-flow date. It represents the date on which the cumulative streamflow reaches half of the annual discharge. |
| hfd_std | day of year | Standard deviation of the mean half-flow dates. |
| q_5 | mm day ⁻¹ | 5% flow quantile, which represents low flows. |
| q_95 | mm day ⁻¹ | 95% flow quantile, which represents high flows. |
| hq_freq | days yr ⁻¹ | Frequency of Q > 9 times the median daily flow. |
| hq_dur | days | Average duration of flow events of consecutive days >9 times the median daily flow. |
| lq_freq | days yr ⁻¹ | Frequency of Q < 0.2 times the median daily flow. |
| lq_dur | days | Average duration of flow events of consecutive days <0.2 times the median daily flow. |
| zero_q_freq | — | Frequency of days with Q = 0 |
| p_mean | mm day ⁻¹ | Mean daily precipitation. |
| pet_mean | mm day ⁻¹ | Mean daily potential evapotranspiration (PET). |
| aridity | — | Ratio between PET and precipitation. |
| p_seasonality | — | Seasonality and timing of precipitation, which was estimated using the precipitation and temperature time series, and computed using Eq. (13) in Woods, 2009 ¹⁰² . |
| frac_snow | — | Fraction of precipitation falling as on days colder than 0°C. |
| hp_freq | days yr ⁻¹ | Frequency of P > 5 times the median daily precipitation (high precipitation events). |
| hp_dur | days | Average duration of periods with consecutive high precipitation events. |
| hp_time | season | Season during most high precipitation events occur (e.g., Fall, Winter, Summer or Spring). |
| lp_freq | days yr ⁻¹ | Frequency of P events < 1 mm day ⁻¹ (dry days). |
| lp_dur | days | Average duration of periods with consecutive dry days. |
| lp_time | season | Season during most dry days occur (e.g., Fall, Winter, Summer or Spring). |
| num_years_{hydro, climatic} | — | Number of years with hydrological or meteorological observations used for the signatures' computation. |
| start_date_{hydro, climatic} | date | First date with with hydrological or meteorological observations used for the signatures' computation. |
| end_date_{hydro, climatic} | date | Last date with hydrological or meteorological used for the signatures' computation. |

Table 5. Set of static hydro-climatic signatures. The hydrological year considered in this study starts at 1st of October and goes until the 30th of September. Unlike streamflow indices, these signatures are static, each represented by a single value calculated for the available data for the period from 1950 to 2022.

| Group | Attribute | Description | Unit | Source |
|---|--|--|----------------------|---------------------|
| Meteorology | p_mean | Total mean daily precipitation measured as the height of the equivalent liquid water in a square meter. | mm day ⁻¹ | E-OBS ¹⁸ |
| | t_{mean, min, max} | Daily mean, minimum and maximum air temperature measured near the surface. | °C | |
| | sp_mean | Mean air pressure at sea level. | hPa | |
| | rh_mean | Daily mean relative humidity measured near the surface. | % | |
| | ws_mean | Daily mean wind speed at 10-meter height. | ms ⁻¹ | |
| | swr_mean | The flux of shortwave radiation (also known as solar radiation) measured at the Earth's surface. | Wm ⁻² | |
| | pet_mean | Potential evapotranspiration was estimated using the Hargreaves equation ⁴³ . | mm day ⁻¹ | derived |
| | stations_num_{p_mean, t_mean, t_min, t_max, sp_mean, rh_mean, ws_mean, swr_mean} | Number of weather stations measuring the given variable within the catchment boundary assuming a 10 km buffer. | - | E-OBS ¹⁸ |
| stations_dens_{p_mean, t_mean, t_min, t_max, sp_mean, rh_mean, ws_mean, swr_mean} | Weather stations density for the given variable within the catchment boundary. | Stations km ⁻² | | |

Table 6. Meteorological catchment attributes at daily resolution from 1950 to 2022. These attributes are aggregated over individual catchment boundaries. The table details both the time series variables and the information regarding the number of stations and their density.

| Group | Attribute | Description | Unit | Source |
|------------|-----------------------------------|--|--------------------------------|--|
| Topography | ele_mt_{max, mean, min} | Mean, minimum and maximum elevation. | m | MERIT-Hydro ^{21,24} |
| | slp_dg_mean | Mean terrain slope. | ° | |
| | flat_area_fra | Percentage of area with slope <3°. | % | |
| | steep_area_fra | Percentage of area with slope >15°. | % | |
| | elon_ratio | Derived elongation ratio ¹⁰³ | — | |
| | strm_dens | Stream density, ratio of lengths of streams and the catchment area. | 1000 Km km ⁻² | |
| Soils* | root_dep | Depth available for roots. | cm | European Soil Database Derived data (ESDD) ²⁸⁻³⁰ |
| | soil_tawc | Total available water content. | mm | |
| | soil_fra_{sand, silt, clay, grav} | Sand, silt, clay and gravel fraction of soil material. | % | |
| | soil_bd | Bulk density. | g cm ⁻³ | |
| | oc_fra | Fraction of organic material. | % | |
| Geology | lit_fra_{class} | Percentage of each lithological class aggregated over the catchment. | % | Global Lithological Map Database (GLiM) ¹⁹ |
| | lit_dom | Lithological dominant class. | Classes (n = 16) | |
| | tot_area | Percentage of the catchment area covered by GLiM. | % | |
| | bedrk_dep | Depth to bedrock. | m | Pelletier, 2016 ²⁰ |
| Hydrology | dam_num | Number of dams upstream. | — | Georeferenced global Dams and Reservoirs ²² |
| | res_num | Number of reservoirs upstream. | — | |
| | dam_yr_{first, last} | First and last years of dam's construction. | — | |
| | res_tot_sto | Total upstream storage volume. | 10 ⁶ m ³ | |
| | lakes_num | Number of lakes upstream. | — | HydroLakes ⁴⁵ |
| | lakes_tot_area | Total area covered by lakes upstream. | Km ² | |
| | lakes_tot_vol | Total upstream volume. | 10 ⁶ m ³ | |
| Vegetation | ndvi_{month, mean}** | Mean NDVI over the catchment area. | — | MODIS ³¹ |
| | lai_{month, mean}** | Mean LAI over the catchment area. | — | MODIS ³² |
| Landcover | sno_cov_{month, mean}** | Mean snow cover percentage over the catchment area. | % | MODIS ²⁷ |

Table 7. Set of static catchment attributes included in the present dataset. *All soil attributes were aggregated by mean, max, min, P05, P25, med, P75 and P90, which sums to a total of 64 variables. **NDVI, LAI and snow cover attributes were aggregated considering the total mean and the month of the year (January = 01 to December = 12) mean from the period between 01.01.2001 to 31.12.2022, which means that each attribute has 13 variables here referred as static since not shown in a time series format.

- **meteorology:** Contains one csv-file per catchment (17,130 in total), each containing all the daily aggregated meteorological forcing records for that catchment (as detailed in Table 6). The rows of each csv-file represent the time, and the columns represent each of the 9 meteorological variables.
- **attributes:** Contains two sub folders. The **static_attributes** subfolder contains one csv-file per attribute group (i.e., topography, soils, geology, hydrology, vegetation and landcover) encompassing all the attributes shown in Table 7. The rows of the csv-file represent the gauging stations, and the columns represent the attribute variable. The **temporal_attributes** subfolder includes all the monthly or annual landscape attributes shown in Table 8. The csv-files in this subfolder are organized by gauging stations (rows), and attribute variables (columns), or as time series (each column represents one gauging station, and each row represents one date).
- **hydroclimatic_signatures:** Contains one csv-file with all computed hydro-climatic signatures for all catchments. The rows of each csv-file represent the streamflow gauging station, and the columns represent each of the 25 derived signatures.
- **appendix:** Contains three txt-files. One file provides descriptions of the lithological classes' labels, another describes the landcover classes' labels, and the third file includes licenses and data providers.

Streamflow data catalogue. An important component of EStreams is the streamflow catalogue, which provides complete guidance on how to retrieve the raw streamflow data used in this study to compute the streamflow statistics. Table 9 provides an overview and description of the attribute fields included in the catalogue.

Particularly, the field **license_redistribution** specifies the data redistribution policy of the data provider. In cases where this information is unavailable, users are advised to proceed with caution regarding any redistribution or specific use of the data, and to contact the data provider directly. The catalogue also includes various links to individual data providers, covering the website, the license source, streamflow and gauges metadata. Up to four different links are provided because the websites for downloading the streamflow time series may differ from those for the gauges metadata.

The Zenodo repository⁴⁶ (<https://doi.org/10.5281/zenodo.13154470>) supports versioning, which ensures reproducibility, benchmarking, and the extensibility of the dataset as new stations or time periods are added.

Additionally, Jupyter Notebook demonstrations are available at the GitHub repository⁴⁷ (<https://doi.org/10.5281/zenodo.13255133>) showing not only how to use the catalogue but also allowing to directly retrieve

| Group | Attribute | Description | Unit | Source |
|------------|--------------------|---|------------------|---|
| Vegetation | ndvi_mean | Monthly and yearly NDVI. | — | MODIS ³¹ |
| | lai_mean | Monthly and yearly LAI. | — | MODIS ³² |
| Landcover | sno_cov_mean | Monthly and yearly snow cover percentage time series. | % | MODIS ²⁷ |
| | irrig_area_{yr} | 10/5-year resolution total area equipped for irrigation. | km ² | AEI_EARTHSTAT_IR product from HID ²⁶ |
| | tot_area_{year} | Fraction of the catchment area covered by the Corine product. | — | CORINE ²⁵ |
| | luc_dom_{year} | Land cover majority class for 1990, 2000, 2006, 2012 and 2018. | Classes (n = 44) | |
| | luc_{year}_{class} | Fraction of each landcover class aggregated over the catchment for 1990, 2000, 2006, 2012 and 2018. | — | |

Table 8. Set of the temporal catchment landscape attributes. Vegetation and snow cover attributes have a monthly and yearly resolution from 2001–2022. The irrigation has a variable window resolution of 10–5-years from 1900–2005.

| Attribute name | Description |
|------------------------|---|
| provider_id | Unique code used to refer the <i>basin_id</i> to their respective data provider |
| code_basins | Code shown in the first two-four digits of the <i>basin_id</i> of their respective catchments |
| provider_country | Country name of the data provided. |
| country_code | Country code of the data provided (e.g., PT for Portugal or AT for Austria). |
| provider_name | Name of the data provider. |
| license_redistribution | Type of redistribution license. |
| platform | Platform where the dataset is available. Either a website, or via contact request. |
| num_stations | Total number of streamflow stations available on the platform as of the date the catalogue data was derived. |
| start_date | Date of the first available streamflow measurement at the date of request/download. |
| end_date | Date of the last available streamflow measurement at the date of request/download. |
| website | Link to the official website of the data provider. |
| source_license | Link where the users can get further information regarding license and terms of use (when available). |
| source_streamflow | Link to the streamflow data provider website. |
| source_gauges_infos | Link to the official source where the gauges information is available (location, river and name). |
| references | Formal reference for citing the streamflow data. |
| observations | Extra information when needed to provide further guidance to the users. |
| download_method | Method of download available at the moment of publication. This specifies if users should download the data manually and individually, or if there is an official API, a provided code, or if a contact form is necessary to request the records. |

Table 9. Attribute fields included in the European Streamflow Catalogue provided.

and pre-process each of the daily records currently included in EStreams. The repository is linked to a GitHub page, enabling users to track potential changes in data providers, websites, and propose updates. This collaborative approach can lead to new releases of the catalogue, ensuring EStreams remains an updated and dynamic resource.

Gauges layer. A comprehensive overview of the gauges' attributes and metadata included in this dataset is presented in Table 10. These attributes are designed to offer users complete guidance on data availability before downloading, thereby optimizing the data collection process. The attributes include the gauges names and location, data provider, topographic information, temporal data availability, quality and reliability descriptors, and nested catchments & flow order attributes. These attributes ensure that users have detailed information to facilitate the efficient retrieval and application of the streamflow data in various hydrological analyses.

Catchments layer. The delineated boundary of each catchment is stored in the catchment layer. This layer includes the *basin_id* field, which is also used for the gauges, allowing a link between the two datasets. Additionally, the catchment layer also has the fields *gauge_id*, *gauge_country* (here named *country*), *area_official* (here named *area_offic*), *area_estreams* (here named *area_estre*), *area_flag*, *area_rel*, *start_date*, *end_date*, *gauge_flag*, *gauges_upstream* (here named *upstream*) and *watershed_group* (here named *group*), which were already described in Table 10. Note that *area_official*, *area_estreams*, *gauge_country*, *gauges_upstream* and *watershed_group* had their names reduced due to storage limitations in the shape files. These fields ensure consistency between the catchment and gauge datasets, facilitating seamless integration and analysis.

Technical Validation

Duplicate stations. This work provides, alongside the gauges' metadata, information on potential candidates for duplication. This information is useful for users aiming to have a consistent dataset for their hydrological analysis. The results indicate that a total of 885 gauges are identified as potential duplicates, representing about 5%

| Attribute name | Description |
|---------------------|---|
| basin_id | An 8-digit code defined by this work. |
| gauge_id | The official code available by the data source, which can be used to retrieve records directly from the data providers. |
| gauge_name | The official name of the station provided by the data source*. |
| gauge_country | Country code where the gauge is located (e.g., PT for Portugal or AT for Austria). |
| gauge_provider | Data source code aligned with the catalogue. |
| river | The name of the river provided by the data source*. |
| lon_snap | Longitude of the gauge in WGS84 original or moved. |
| lat_snap | Latitude of the gauge in WGS84 original or moved. |
| lon | Longitude of the gauge in WGS84 provided by the data source. |
| lat | Latitude of the gauge in WGS84 provided by the data source. |
| elevation | The official gauge elevation reported by the data provider*. |
| area_official | The official area reported by the data provider (A_{official})*. |
| area_estreams | The area (in km ²) derived from the current delineation methodology (A_{EStreams}). |
| area_flag | A quality flag for the current area computation as reported in Table 3. |
| area_rel | The percentual (%) relative difference between the derived and the reported area, relative to the reported area, as defined by Eq. (1). |
| start_date | First date with valid observations as of the date the data was accessed. |
| end_date | Last date with valid observations as of the date the data was accessed. |
| num_years | Number of years with valid data. |
| num_months | Number of months with valid data. |
| num_days | Number of days with valid data. |
| num_continuous_days | Maximum number of days between the <i>start_date</i> and <i>end_date</i> with no gaps. |
| num_days_gaps | Number of days with gaps between the <i>start_date</i> and <i>end_date</i> . |
| num_days_reliable | Number of days with data classified as “reliable” from the respective provider. |
| num_days_noflag | Number of days with data without a quality flag provided by the respective provider. |
| num_days_suspect | Number of days with data classified as “suspect” from the respective provider. |
| gauge_flag | Quality flag of the respective streamflow gauge as reported in Table 2. |
| duplicated_suspect | If it is the case, <i>basin_id</i> of the gauge suspect of being a duplicate with this gauge. |
| watershed_group | A number assigning to which main watershed is the gauge belongs to, e.g., all gauges within the Rhine watershed are assigned the number 1. |
| gauges_upstream | The number of unique gauging stations upstream of the given gauge. This count includes the basin itself but excludes any duplicate stations. This means that if one gauge has a duplicate, the count considers only one gauge. |
| nested_catchments | A list of all nested catchments within the given basin. This list includes the basin itself and may differ from the total number in <i>gauges_upstream</i> because it includes all gauges, retaining any duplicates within the same list. |

Table 10. Description of the attributes of the streamflow gauges’ layer. *These are information seldom not available from official sources.

of the total. This means that more than 16,600 gauges in the dataset may be seen as unique gauging stations. The duplicates are divided into two types: gauges duplicated with other gauges within the same provider and gauges duplicated with other gauges within different providers.

These first types of duplicates often occur when gauges are discontinued and later reactivated as new stations, usually resulting in stations with non-overlapping time records but located at the same point. These cases are primarily found in France (449) and Finland (160). For example, stations FR001479 (1969–1999), FR001477 (1993–1999) and FR001478 (2015–2023) are flagged as duplicate suspects among each other.

Additionally, 163 gauges are identified as duplicates across different data providers. These typically represent gauging stations located at the boundaries between countries and are mainly found in Austria (33), Switzerland (36) and Czech Republic (51). Interestingly, FR004543 is the only gauge identified as duplicate both within the same provider (FR002217) and across different providers (CH000268).

Basin delineation validation. In this part of the study, we used the dataset provided by LamaH-CE⁸ for Austria, which includes both catchment boundaries and their respective officially reported areas. These were compared to the boundaries delineated using the methodology adopted in this work.

Figure 5a shows a scatter plot comparing the areas reported in LamaH-CE and those derived in EStreams. As expected, the scatter between the computed and reported areas is larger for smaller catchments. Figure 5b presents a histogram with the distribution of the relative absolute area difference $|A_{\text{rel}}|$ between the two areas (in %). Out of the total of 599 Austrian catchments, 539 had a $|A_{\text{rel}}|$ below 10%. This indicates that roughly 90% of the catchments were accurately delineated during the automatic part of the delineation process.

However, if we consider only catchments with areas above 100 km² the number of catchments with $|A_{\text{rel}}|$ above 10% drops from 60 to only 21. After visual inspection, we concluded that the main cause of these discrepancies was associated either to the difficulties in the delineation of relatively small catchments, below 100 km², or to small discrepancies between the streamflow gauge location in terms of the MERIT-Hydro network.

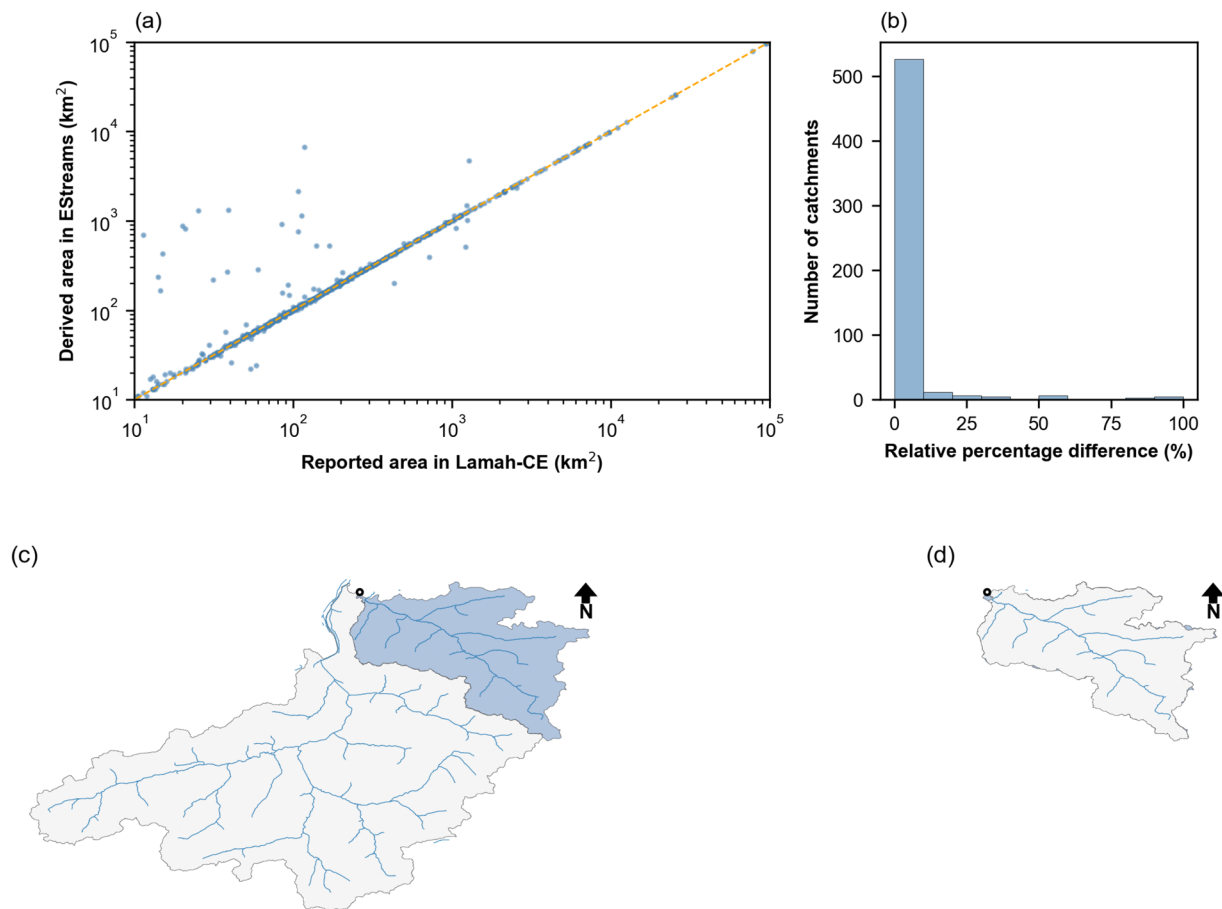


Fig. 5 (a) Comparison of catchment boundary areas reported LamaH-CE⁸ against those delineated in this study. Both axes are presented in logarithmic scale to enhance visualization. (b) Histogram illustrating the $|A_{rel}|$ between the two sources of data. Most catchments exhibit $|A_{rel}|$ below 10%. Catchment AT000009 (EStreams) delineations are displayed (c) prior to manual adjustment of the outlet location and (d) following manual adjustment.

Figure 5c-d illustrate an example of the catchment delineation workflow for catchment AT000009. This catchment has an $A_{official}$ of 1281.0 km². Initially, $A_{EStream}$ derived an area of 4680.0 km², which accounts for a A_{rel} of +265.0%. Upon visual inspection, we realized that the inconsistency was due to the inaccurate location of the streamflow gauge in relation to the MERIT-Hydro River network (Fig. 5c). Since the outlet was not within the river network, the “delineator” python module used automatically moved it to the closest river network intersection, which had a much higher drainage area. After manually adjusting the streamflow gauge location, the delineation resulted in an area of 1,300.0 km², an A_{rel} of only +1.5% (Fig. 5d).

E-OBS assessment. Spatial coverage. EStreams used E-OBS to derive the catchment aggregated time series of meteorological variables. However, the number of stations used to produce the gridded dataset varies significantly from country to country. Here we provide a brief overview of the station densities used to derive the precipitation time series provided in E-OBS within each catchment. We present this analysis only for precipitation since it is considered the most important forcing input in hydrological studies and gives already a significant overview of the E-OBS network. To ensure a fair comparison, we considered a buffer of 10 km for the catchment boundaries and considered any station within this range to compute the number of stations.

Figure 6a illustrates the spatial distribution of the stations, revealing a large spatial variability in station density. Central and North Europe exhibit the highest density, with Germany and Poland taking leading in station density, while the density decreases significantly towards South and East.

Figure 6b presents the histogram of the station density per catchment included in EStreams. The x-axis is resampled to stations per 100 km² to facilitate visualization, with the threshold of less than one station per 100 km² marked in red. A total of 9,840 catchments have at least one precipitation gauge per 100 km². This represents, a median of 1.2 stations per 100 km². Considering absolute terms, we found a total of 14,153 gauges with at least one precipitation station within their boundaries.

This information enables users to be aware of the highly variable quality of the provided E-OBS data and make informed decisions, especially considering the critical role of accurate precipitation data in many hydrological applications. Like streamflow data, national providers typically offer much higher resolution precipitation

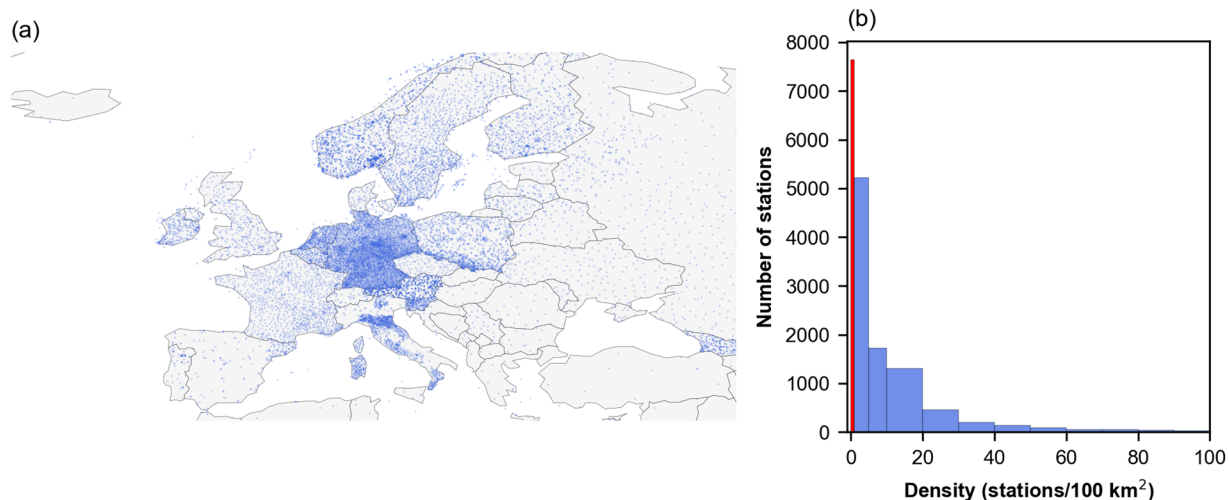


Fig. 6 (a) Overview of the spatial distribution of the stations used to derive the precipitation time series gridded data available at E-OBS¹⁸. (b) Histogram of the stations per catchment. Due to the high distribution of densities the bins are not evenly spaced, and the first bin (in red) corresponds to the threshold of one station per 100 km². Basemap from GeoPandas¹⁰⁴.

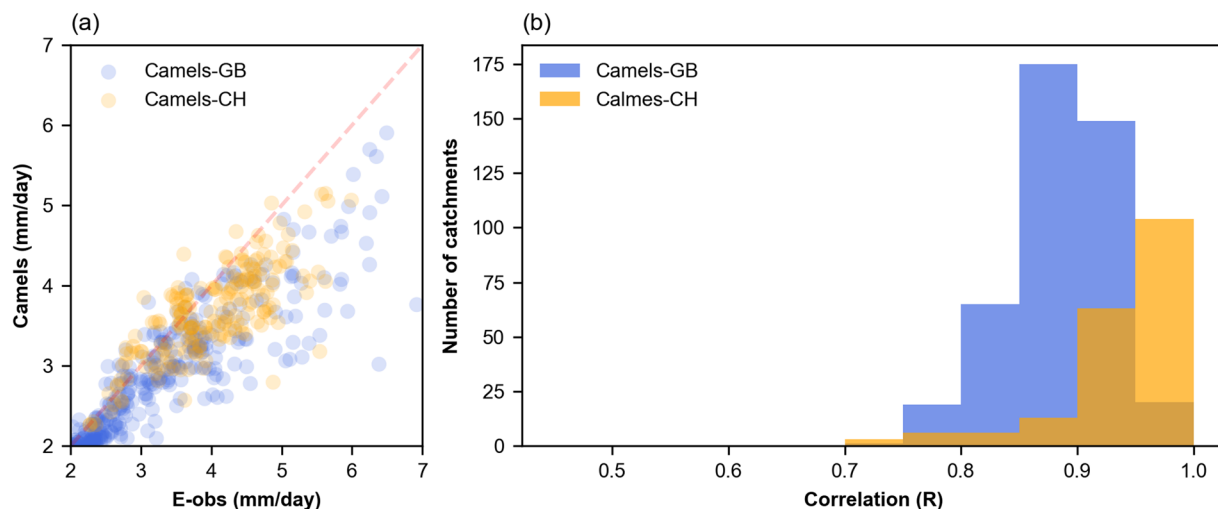


Fig. 7 (a) Scatter plot of the long-term mean daily precipitation (1950–2022) considering the precipitation forcing time series derived from E-OBS¹⁸ and the provided in CAMELS-CH⁷ and CAMELS-GB² and (b) Histogram of the correlation coefficient between the two data sources. The plots only show catchments with areas above 100 km².

data compared to global databases⁴⁸. While retrieving this information was beyond the scope of this study, users may choose to leverage such local data sources, particularly in regions where station density is notably low, such as in the South, East, and West of Europe.

Validation of meteorological forcing. We further validated the aggregated precipitation derived from E-OBS comparing it to the reported time series available at CAMELS-CH⁷ and CAMELS-GB². Given that the aggregation of the forcing variables used E-OBS gridded data with a resolution of 0.25 degrees, we opted to include only catchments with areas above 100 km² in the comparison.

Figure 7a shows a scatter plot illustrating the daily precipitation from E-OBS and CAMELS. CAMELS-GB is represented in blue and CAMELS-CH in orange. A notable correspondence between the two sources is observable, with correlation coefficients of 0.89 for GB and 0.94 for CH. Generally, the scatter is lower in catchments with higher daily mean precipitation and an underestimation from E-OBS compared to the two sources is evident.

Figure 7b shows the distribution of the correlation coefficients between each daily time series of E-OBS and CAMELS. Again, it is possible to observe that most of the catchments presented a correlation above 0.8, indicating some agreement between the two precipitation sources. Overall, CAMELS-CH demonstrates higher correlation coefficients than CAMELS-GB. Despite this comparison only encompassing two different regions

within the large span covered by EStreams, it was conducted using two independent sources. Hence, this analysis suggests that E-OBS, at least in countries where the station density is relatively high, provides a broadly consistent starting point for representing precipitation time series.

Usage Notes

Aggregated data. The original data used to aggregate the catchment attributes such as climate, geology, hydrology, land use and land cover, soil types and vegetation characteristics have all continental or global resolution. It should be kept in mind that such resolution is rather coarse compared to local information usually available at the national scales, but seldom easily accessible. We therefore recommend that users acknowledge these potential limitations when using the aggregated data. Additionally, we recommend users to also reference the original sources when using the aggregated data provided in EStreams.

Streamflow catalogue. We recognize that potential retrospective check and updates of streamflow time series by the data providers may alter the information of the gauges metadata provided here. We also acknowledge that potential changes in the data providers' platforms may alter the available links in the catalogue. Therefore, we invite the users to access the latest version of the catalogue and dataset on the Zenodo repository⁴⁶ page for potential updates.

Instructions for Python. We kindly request that future users of the EStreams' codes read and follow carefully the instructions provided in the scripts. Specifically, (i) use the specified version of the Python modules (requirements.txt); (ii) clone the repository locally and keep all the original folders' names; (iii) place the original data in their specified folder and with their expected filename and version; (iv) follow the pre-defined specified order of run for the available scripts (when necessary). Be aware that the potential main source of problems when running the scripts might be caused by not following these guidelines.

Code availability

The current version of the code used to produce the EStreams dataset and catalogue (v1.0.0) is available at a Zenodo repository⁴⁷ at <https://doi.org/10.5281/zenodo.13255133>. For the latest version of the code, users are invited to visit the project GitHub repository at <https://github.com/thiagovmdon/EStreams>. The scripts are organized to enable users to follow a logical sequence during code usage. All data processing scripts are written in Python, while some data retrieval tasks are performed using JavaScript for the Google Earth Engine (GEE) platform. Although all scripts are executable, users must download and preprocess the original data due to redistribution licenses. Detailed instructions regarding the version used, data retrieval, and any required preprocessing are provided within the respective scripts.

Received: 4 March 2024; Accepted: 29 July 2024;

Published online: 13 August 2024

References

1. Addor, N., Newman, A. J., Mizukami, N. & Clark, M. P. The CAMELS data set: Catchment attributes and meteorology for large-sample studies. *Hydrol Earth Syst Sci* **21**, 5293–5313 (2017).
2. Coxon, G. *et al.* CAMELS-GB: hydrometeorological time series and landscape attributes for 671 catchments in Great Britain. *Earth Syst Sci Data* **12**, 2459–2483 (2020).
3. Kratzert, F. *et al.* Caravan - A global community dataset for large-sample hydrology. *Scientific Data* **10**, 1–11 (2023).
4. Fowler, K. J. A., Acharya, S. C., Addor, N., Chou, C. & Peel, M. C. CAMELS-AUS: Hydrometeorological time series and landscape attributes for 222 catchments in Australia. *Earth Syst. Sci Data* **13**, 3847–3867 (2021).
5. Chagas, V. B. P. *et al.* CAMELS-BR: Hydrometeorological time series and landscape attributes for 897 catchments in Brazil. *Earth Syst Sci Data* **12**, 2075–2096 (2020).
6. Alvarez-Garreton, C. *et al.* The CAMELS-CL dataset: Catchment attributes and meteorology for large sample studies-Chile dataset. *Hydrol Earth Syst Sci* **22**, 5817–5846 (2018).
7. Höge, M. *et al.* CAMELS-CH: hydro-meteorological time series and landscape attributes for 331 catchments in hydrologic Switzerland. *Earth Syst Sci Data* **15**, 5755–5784 (2023).
8. Klingler, C., Schulz, K. & Herrnegger, M. LamaH-CE: LARge-SaMple DAta for Hydrology and Environmental Sciences for Central Europe. *Earth Syst Sci Data* **13**, 4529–4565 (2021).
9. Arsenaault, R. *et al.* A comprehensive, multisource database for hydrometeorological modeling of 14,425 North American watersheds. *Scientific Data* **7**, 1–12 (2020).
10. Hao, Z. *et al.* CCAM: China Catchment Attributes and Meteorology dataset. *Earth Syst Sci Data* **13**, 5591–5616 (2021).
11. Marti, B. *et al.* CA-discharge: Geo-Located Discharge Time Series for Mountainous Rivers in Central Asia. *Scientific Data* **10**, 1–21 (2023).
12. Helgason, H. B. & Nijssen, B. LamaH-Ice: LARge-SaMple DAta for Hydrology and Environmental Sciences for Iceland, CUAHSI HydroShare, <https://www.hydroshare.org/resource/86117a5f36cc4b7c90a5d54e18161c91/> (last access: 01 May) (2024).
13. Do, H. X., Gudmundsson, L., Leonard, M. & Westra, S. The Global Streamflow Indices and Metadata Archive (GSIM)-Part 1: The production of a daily streamflow archive and metadata. *Earth Syst Sci Data* **10**, 765–785 (2018).
14. Gudmundsson, L., Do, H. X., Leonard, M. & Westra, S. The Global Streamflow Indices and Metadata Archive (GSIM)-Part 2: Quality control, time-series indices and homogeneity assessment. *Earth Syst Sci Data* **10**, 787–804 (2018).
15. Chen, X., Jiang, L., Luo, Y. & Liu, J. A global streamflow indices time series dataset for large-sample hydrological analyses on streamflow regime (until 2022). *Earth Syst Sci Data* **15**, 4463–4479 (2023).
16. GRDC. Global Runoff Data Center: River discharge data. Federal Institute of Hydrology, 56068 Koblenz, Germany. <https://www.bafg.de/GRDC> (last access: 01 May 2024).
17. Färber, C. *et al.* GRDC-Caravan: extending the original dataset with data from the Global Runoff Data Centre (0.1) [Data set]. Zenodo <https://zenodo.org/records/8425587>, <https://doi.org/10.5281/ZENODO.8425587> (2023).
18. Cornes, R. C., van der Schrier, G., van den Besselaar, E. J. M. & Jones, P. D. An Ensemble Version of the E-OBS Temperature and Precipitation Data Sets. *Journal of Geophysical Research: Atmospheres* **123**, 9391–9409 (2018).
19. Hartmann, J., Moosdorf, N., Hartmann, J. & Moosdorf, N. The new global lithological map database GLiM: A representation of rock properties at the Earth surface. *Geochemistry, Geophysics, Geosystems* **13**, 12004 (2012).

20. Pelletier, J. D. *et al.* A gridded global data set of soil, intact regolith, and sedimentary deposit thicknesses for regional and global land surface modeling. *J Adv Model Earth Syst* **8**, 41–65 (2016).
21. Yamazaki, D. *et al.* MERIT Hydro: A High-Resolution Global Hydrography Map Based on Latest Topography Dataset. *Water Resour Res* **55**, 5053–5073 (2019).
22. Wang, J. *et al.* GeoDAR: georeferenced global dams and reservoirs dataset for bridging attributes and geolocations. *Earth Syst Sci Data* **14**, 1869–1899 (2022).
23. Linke, S. *et al.* Global hydro-environmental sub-basin and river reach characteristics at high spatial resolution. *Scientific Data* **6**, 1–15 (2019).
24. Yamazaki, D. *et al.* A high-accuracy map of global terrain elevations. *Geophys Res Lett* **44**, 5844–5853 (2017).
25. CORINE: CORINE Land Cover — Copernicus Land Monitoring Service. European Environment Agency [data set], Copenhagen, Denmark, <https://land.copernicus.eu/en/products/corine-land-cover>.
26. Siebert, S. *et al.* A global data set of the extent of irrigated land from 1900 to 2005. *Hydrol Earth Syst Sci* **19**, 1521–1545 (2015).
27. Hall, D. K. & Riggs, G. A. MODIS/Terra Snow Cover Daily L3 Global 500m SIN Grid, Version 61 [Data Set]. NASA National Snow and Ice Data Center Distributed Active Archive Center. vol. 21. <https://doi.org/10.5067/MODIS/MOD10A1.061> (2021).
28. Hiederer, R. *Mapping Soil Typologies – Spatial Decision Support Applied to European Soil Database*. <https://doi.org/10.2788/87286> (2013).
29. Hiederer, R. *Mapping Soil Properties for Europe – Spatial Representation of Soil Database Attributes*. <https://data.europa.eu/doi/10.2788/94128> (2013).
30. ESDD. *European Soil Database Derived Data*. <https://esdac.jrc.ec.europa.eu/Content/European-Soil-Database-Derived-Data> (Last Access: 23 Nov 2023).
31. Didan, K. MODIS/Terra Vegetation Indices 16-Day L3 Global 500m SIN Grid V061 [Data set]. ASA EOSDIS Land Processes Distributed Active Archive Center <https://doi.org/10.5067/MODIS/MOD13A1.061> (2021).
32. Myneni, R., Knyazikhin, Y. & Park, T. MODIS/Terra Leaf Area Index/FPAR 8-Day L4 Global 500m SIN Grid V061 [Data set]. NASA EOSDIS Land Processes Distributed Active Archive Center, <https://doi.org/10.5067/MODIS/MOD15A2H.061> (2021).
33. Christen, P. Data matching: Concepts and techniques for record linkage, entity resolution, and duplicate detection. *Data Matching: Concepts and Techniques for Record Linkage, Entity Resolution, and Duplicate Detection* 1–270. <https://doi.org/10.1007/978-3-642-31164-2/COVER> (2012).
34. Trambly, Y. *et al.* ADHI: The African Database of Hydrometric Indices (1950–2018). *Earth Syst Sci Data* **13**, 1547–1560 (2021).
35. Crochemore, L. *et al.* Lessons learnt from checking the quality of openly accessible river flow data worldwide. *Hydrological Sciences Journal* **65**, 699–711 (2020).
36. Heberger, M. *delineator.py: fast, accurate global watershed delineation using hybrid vector- and raster-based methods*. Zenodo <https://doi.org/10.5281/ZENODO.7314287> (2022).
37. COPERNICUS Land Monitoring Service. EU-Hydro. <https://land.copernicus.eu/imagery-in-situ/eu-hydro> (last access: 18 Aug 2023). 2019.
38. Dal Molin, M. dalmo1991/HydroAnalysis: v1.0.0 (1.0.0). Zenodo <https://doi.org/10.5281/zenodo.5716016> (2021).
39. Wunsch, A. *et al.* Karst spring discharge modeling based on deep learning using spatially distributed input data. *Hydrol Earth Syst Sci* **26**, 2405–2430 (2022).
40. Rojas, R., Feyen, L., Dosio, A. & Bavera, D. Improving pan-European hydrological simulation of extreme events through statistical bias correction of RCM-driven climate simulations. *Hydrol Earth Syst Sci* **15**, 2599–2620 (2011).
41. Becker, A. *et al.* A description of the global land-surface precipitation data products of the Global Precipitation Climatology Centre with sample applications including centennial (trend) analysis from 1901–present. *Earth Syst Sci Data* **5**, 71–99 (2013).
42. Bandhauer, M. *et al.* Evaluation of daily precipitation analyses in E-OBS (v19.0e) and ERA5 by comparison to regional high-resolution datasets in European regions. *International Journal of Climatology* **42**, 727–747 (2022).
43. Hargreaves, G. H. & Samani, Z. A. Estimating potential evapotranspiration. *Journal of Irrigation and Drainage Engineering* **108**, 223–230 (1982).
44. Vremec, M. & Collenteur, R. PyEt—a Python package to estimate potential and reference evapotranspiration 1.1.0. in *EGU General Assembly Conference Abstracts* (2021).
45. Messenger, M. L., Lehner, B., Grill, G., Nedeva, I. & Schmitt, O. Estimating the volume and age of water stored in global lakes using a geo-statistical approach. *Nature Communications* **7**, 1–11 (2016).
46. do Nascimento, T. V. M. *et al.* EStreams: An Integrated Dataset and Catalogue of Streamflow, Hydro-Climatic Variables and Landscape Descriptors for Europe (1.0) [Data set], Zenodo, <https://doi.org/10.5281/zenodo.13154470> (2024).
47. do Nascimento, T. V. M. *et al.* EStreams: An Integrated Dataset and Catalogue of Streamflow, Hydro-Climatic Variables and Landscape Descriptors for Europe (v.1.0.0) [Code], Zenodo, <https://doi.org/10.5281/zenodo.13255133> (2024).
48. Clerc-Schwarzenbach, F. M. *et al.* HESS Opinions: A few camels or a whole caravan? *EGU sphere [preprint]* <https://doi.org/10.5194/egusphere-2024-864> (2024).
49. BML. Federal Ministry of Agriculture, Forestry, Regions and Water Management: WebGIS-Applikation eHYD, Wien, Austria, <https://ehyd.gv.at> (last access: 05 May 2023).
50. FHMZBIH. Federalni hidrometeorološki zavod: Početna: idrologija: hidrološki godišnjaci, Bosnia. <https://www.fhmzbih.gov.ba/latinica/HIDRO/godisnjaci.php> (last access: 29 June 2023).
51. VV. Vlaanderen waterinfo, Belgium. <https://www.waterinfo.be/kaartencatalogus?KL=en> (last access: 07 Dec 2023).
52. SPW. Service public de Wallonie: L'hydrométrie en Wallonie: Observations: Debit, Belgium. <https://hydrometrie.wallonie.be/home/observations/debit.html?mode=announcement> (last access: 07 Dec 2023).
53. BAFU. Federal Office for the Environment, Switzerland. <https://www.bafu.admin.ch/bafu/en/home.html> (last access: 15 May 2023).
54. CHMI. Czech Hydrometeorological Institute: ISVS - Evidence množství povrchových vod. [https://isvs.chmi.cz/ords/f?p=11002:HOME:5026647009329:::~](https://isvs.chmi.cz/ords/f?p=11002:HOME:5026647009329:::) (last access: 10 Jul 2023).
55. LHW. Landesbetrieb für Hochwasserschutz und Wasserwirtschaft Sachsen-Anhalt, <https://gld.lhw-sachsen-anhalt.de/> (last access: 12 Dec 2023).
56. ASOEAG. Saxon State Office for Environment, Agriculture and Geology: Datenportal für Umweltdaten Sachsen (iDA), https://www.umwelt.sachsen.de/umwelt/infosysteme/ida/processingChain?conditionValuesSetHash=0A8BBED&selector=ROOT.Thema%20Wasser.Oberirdische%20Gew%C3%A4sser.Pegel.Wasserstand%20und%20Durchfluss.OWMN%3Aowmn_menge_tagesmittelwerte_v2.sel&sourceOrderAsc=false&columns=9dfa2224-c924-4328-9805-1d34cd748026&offset=0&limit=2147483647&executionConfirmed=false (last access: 12 Dec 2023).
57. Umweltportal. Schleswig-Holstein, Germany. https://umweltportal.schleswig-holstein.de/kartendienste?lang=de&topic=thessd&bgLayer=sgx_geodatenzentrum_de_de_basemapde_web_raster_grau_DE_EPSG_25832_ADV&E=567583.34&N=5998716.15&zoom=4&layers=262b5c716ef5358fc1ac1e34afd45915 (last access: 12 Dec 2023).
58. ELWAS-WEB. Ministerium für Umwelt, Naturschutz und Verkehr des Landes Nordrhein-Westfalen, <https://www.elwasweb.nrw.de/elwas-web/data-objekt;jsessionid=DADDD7196B89E206917D18793294E375;jsessionid=F76CC7CC8ECFBA5F518ECD241AF0BA78?art=Pegel> (last access: 12 Dec 2023).
59. NLWKN. Niedersächsischer Landesbetrieb für Wasserwirtschaft, Küsten- und Naturschutz, <http://www.wasserdaten.niedersachsen.de/cadenza/pages/selector/index.xhtml;jsessionid=1E0F808EF58258C4EE5C777447D1ED4A> (last access: 12 Dec 2023).

60. HLNUG. Hessisches Landesamt für Naturschutz, Umwelt und Geologie. <https://www.hlnug.de/static/pegel/wiskiweb3/webpublic/#/overview/Wasserstand?mode=table&filter=%7B%7D> (last access: 12 Dec 2023).
61. GKD. Bavarian State Office for the Environment – Hydrographic Service, Munich, Germany <https://www.gkd.bayern.de/en/rivers/discharge/tables> (last access: 12 Dec 2023).
62. LUBW. State Agency for the Environment Baden-Württemberg – Hydrographic Service, Karlsruhe, Germany. <https://udo.lubw.baden-wuerttemberg.de/public/> (last access: 12 Dec 2023).
63. WB. Das Wasserportal Berlin: <https://wasserportal.berlin.de/start.php> (last access: 12 Dec 2023).
64. LBAW. Land Brandenburg Auskunftsplattform Wasser. https://apw.brandenburg.de/th=owm_gkp/ (last access: 12 Dec 2023).
65. MKUEM. Ministerium für Klimaschutz, Umwelt, Energie und Mobilität: Rheinland-Pfalz, Germany. <https://wasserportal.rlp-umwelt.de> (data received: 13 Mar 2023).
66. LUBN. Landesamt für Umwelt, Bergbau und Naturschutz. Hochwasser Nachrichten Zentrale: Freistaat Thüringen. <https://hnz.thueringen.de> (data received: 13 Mar 2023).
67. BFG. Bundesanstalt für Gewässerkunde, Germany. https://www.bafg.de/DE/Home/homepage_node.html (data received: 13 Mar 2023).
68. ODA. Overfladevandsdatabasen: Aarhus University, Denmark. <https://odaforalle.au.dk/login.aspx> (last access: 17 Jul 2023).
69. CEDEX. Centro de Estudios y Experimentación de Obras Públicas: Anuario de aforos 2019–2020, Spain. <https://ceh.cedex.es/anuarioaforos/demarcaciones.asp> (last access: 12 Apr 2023).
70. FEI. Finnish Environmental Institute, Finland. <https://www.wp2.ymparisto.fi/scripts/kirjau.asp> (last access: 10 Jul 2023).
71. BanqueHydro. Hydro Portail, France. <https://www.hydro.eaufrance.fr/> (last access: 01 May 2024).
72. NRFA. National River Flow Archive API, United Kingdom. <https://nrfaapps.ceh.ac.uk/nrfa/nrfa-api.html> (last access: 07 Jul 2023).
73. OHIN. Open Hydrosystem Information Network, Greece. <https://openhi.net/en/> (last access: 12 Oct 2023).
74. HCRM. Institute of Marine Biological Resources and Inland Waters, Greece. <https://hydro-stations.hcmr.gr/%cf%80%ce%b1%cf%81%ce%bf%cf%87%ce%ae-%cf%80%ce%bf%cf%84%ce%b1%ce%bc%cf%8e%ce%bd/> (last access: 12 Oct 2023).
75. DHZ. Croatian Meteorological and Hydrological Service. <https://hidro.dhz.hr/> (last access: 01 May 2024).
76. OVF. General Directorate of Water Management. <https://ovf.hu/kozerdeku/adatigenyyles> (data received: 18 Aug 2023).
77. EPA. Environmental Protection Agency, Ireland. <https://epawebapp.epa.ie/hydronet/#Flow> (last access: 27 Jun 2023).
78. OPW. Office of Public Works, Ireland. <https://waterlevel.ie/hydro-data> (last access: 27 Jun 2023).
79. ISPRA. Istituto Superiore per la Protezione e la Ricerca Ambientale, Italy. <http://www.hiscentral.isprambiente.gov.it/hiscentral/hydropmap.aspx?map=obsclient>, (last access: 30 December 2023).
80. APC Abruzzo. Centro Funzionale e Ufficio Idrologia, Idrografico, Mareografico: Agenzia di Protezione Civile della Regione Abruzzo, Italy (data received: 02 August 2023).
81. CFRA Valle d'Aosta. Centro Funzionale Regione Autonoma Valle d'Aosta, Italy. https://presidi2.regione.vda.it/str_dataview_download (last access: 19 May 2023).
82. ARPAE Emilia-Romagna. Agenzia Prevenzione Ambiente Energia - Emilia-Romagna, Italy. <https://simc.arpae.it/dext3r/> (last access: 04 Nov 2023).
83. ARPA Umbria. Agenzia Regionale per la Protezione dell'Ambiente - Umbria, Italy. <https://annali.regione.umbria.it> (last access: 22 May 2023).
84. ARPA Sardegna. Agenzia Regionale per la Protezione dell'Ambiente - Sardegna, Italy. <https://www.sardegnaambiente.it/index.php?xsl=611&s=21&v=9&c=93749&na=1&n=10> (last access: 30 December 2023).
85. ARPA Lombardia. Agenzia Regionale per la Protezione dell'Ambiente - Lombardia, Italy. (data received: 17 Jun 2023).
86. ARPA Lombardia. Agenzia Regionale per la Protezione dell'Ambiente - Lombardia, Italy. https://idro.arpalombardia.it/manual/dati_storici.html (last access: 24 May 2023).
87. ARPA Toscana. Agenzia Regionale per la Protezione dell'Ambiente - Toscana, Italy. <http://www.sir.toscana.it/consistenza-rete> (last access: 16 Jun 2023).
88. ARPA Piemonte. Agenzia Regionale per la Protezione dell'Ambiente - Piemonte, Italy. https://www.arpa.piemonte.it/rischi_naturali/snippets_arpa_graphs/map_meteoweb/?rete=stazione_meteorologica (last access: 22 May 2023).
89. ARPAL Liguria. Agenzia Regionale per la Protezione dell'Ambiente - Liguria, Italy. <https://www.arpal.liguria.it> (data received: 08 Jun 2023).
90. ARPAV Veneto. Agenzia Regionale per la Prevenzione e Protezione Ambientale del Veneto, Italy. <https://www.arpa.veneto.it/> (data received: 30 Jun 2023).
91. SPRUD Trentino. Servizio Prevenzioni Rischi Ufficio Dighe - Trentino-Alto Adige Trento, Italy. <https://www.floods.it/public/DatiStorici.php> (last access: 24 May 2023).
92. NGGL. The National Geoportal of the Grand-Duchy of Luxembourg. <https://map.geoportail.lu> (data received: 13 Mar 2023).
93. RWS. Rijkswaterstaat waterinfo, The Netherlands. <https://waterinfo.rws.nl/#/publiek/waterafvoer> (last access: 07 Dec 2023).
94. NVE. Norwegian Water Resources and Energy Directorate, Norway. <https://seriekart.nve.no> (last access: 10 Jul 2023).
95. IMGW-PIB. Institute of Meteorology and Water Management - National Research Institute, Warszawa, Poland. <https://danepubliczne.imgw.pl/introduction> (last access: 30 Dec 2023).
96. SNIRH. Sistema Nacional de Informação de Recursos Hídricos: Dados de Base, Portugal. <https://snirh.apambiente.pt/index.php?idMain=2&idItem=1> (last access: 01 May 2024).
97. SMHI. Swedish Meteorological and Hydrological Institute, Sweden. <https://www.smhi.se/data/hydrologi/ladda-ner-hydrologiska-observationer#param=waterdischargeDaily,stations=core> (last access: 30 Dec 2023).
98. ARSO. Agencija Republike Slovenije za Okolje, Ljubljana, Slovenia. <https://vode.arso.gov.si/hidarhiv/> (last access: 23 Jun 2023).
99. Sankarasubramanian, A., Vogel, R. M. & Limbrunner, J. F. Climate elasticity of streamflow in the United States. *Water Resour Res* 37, 1771–1781 (2001).
100. Sawicz, K., Wagener, T., Sivapalan, M., Troch, P. A. & Carrillo, G. Catchment classification: Empirical analysis of hydrologic similarity based on catchment function in the eastern USA. *Hydrol Earth Syst Sci* 15, 2895–2911 (2011).
101. Ladson, A. R., Brown, R., Neal, B. & Nathan, R. A standard approach to baseflow separation using the Lyne and Hollick filter. *Australian Journal of Water Resources* 17, 25–34 (2013).
102. Woods, R. A. Analytical model of seasonal climate impacts on snow hydrology: Continuous snowpacks. *Advances in Water Resources* 32(10), 1465–1481, <https://doi.org/10.1016/j.advwatres.2009.06.011> (2009).
103. Schumm, S. A. Evolution of drainage systems and slopes in badlands at Perth Amboy, New Jersey. *GSA Bulletin* 67, 597–646 (1956).
104. Jordahl, K. *et al.* geopandas/geopandas: v0.8.1 *Zenodo* <https://doi.org/10.5281/zenodo.2585848> (2020).

Acknowledgements

This project was funded by a “Money Follows Cooperation” project (Project No. OCENW.M.21.230) between the Netherlands Organization for Scientific Research (NWO) and the Swiss National Science Foundation (SNSF). This work was further supported by the TU Delft Climate Action Research and Education seed funds. We would like to acknowledge all the data providers and contact people who somehow contributed to the construction of this dataset. In particular, we acknowledge specially the E-OBS dataset, the data providers in the ECA&D project (<https://www.ecad.eu>), the European Soil Database Derived data project (ESDAC) and the UK National River Flow Archive.

Author contributions

The co-authors T.N., J.R., R.E., M.H., J.S. M.Hr. and F.F. were involved in the development of the concept of this paper. T.N. and J.R. collected and pre-processed the data. M.C. provided guidance to some data providers in Eastern Europe. T.N. wrote the data aggregation and processing codes in Python and Google Earth Engine. T.N. and J.R. processed the catchment boundaries. T.N. wrote the first draft. M.Hr and F.F. retrieved the funding for the project. All co-authors participated in reviewing the manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41597-024-03706-1>.

Correspondence and requests for materials should be addressed to T.V.M.d.N.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2024