



GRQA: Global River Water Quality Archive

Holger Virro¹, Giuseppe Amatulli^{2,3}, Alexander Knoch¹, Longzhu Shen^{4,5}, and Evelyn Uuemaa¹

¹Department of Geography, Institute of Ecology and Earth Sciences, University of Tartu,
Vanemuise 46, Tartu, 51003, Estonia

²School of the Environment, Yale University, New Haven, CT, 06511, USA

³Center for Research Computing, Yale University, New Haven, CT, 06511, USA

⁴HyperAmp, Barnwell Road, Cambridge CB5 8RQ, UK

⁵Spatial-Ecology, Meaderville House, Wheal Buller, Redruth TR16 6ST, UK

Correspondence: Holger Virro (holger.virro@ut.ee)

Received: 17 February 2021 – Discussion started: 2 March 2021

Revised: 7 September 2021 – Accepted: 18 October 2021 – Published: 30 November 2021

Abstract. Large-scale hydrological studies are often limited by the lack of available observation data with a good spatiotemporal coverage. This has affected the reproducibility of previous studies and the potential improvement of existing hydrological models. In addition to the observation data themselves, insufficient or poor-quality metadata have also discouraged researchers from integrating the already-available datasets. Therefore, improving both the availability and quality of open water quality data would increase the potential to implement predictive modeling on a global scale.

The Global River Water Quality Archive (GRQA) aims to contribute to improving water quality data coverage by aggregating and harmonizing five national, continental and global datasets: CESI (Canadian Environmental Sustainability Indicators program), GEMStat (Global Freshwater Quality Database), GLORICH (GLObal River CHemistry), Waterbase and WQP (Water Quality Portal). The GRQA compilation involved converting observation data from the five sources into a common format and harmonizing the corresponding metadata, flagging outliers, calculating time series characteristics and detecting duplicate observations from sources with a spatial overlap. The final dataset extends the spatial and temporal coverage of previously available water quality data and contains 42 parameters and over 17 million measurements around the globe covering the 1898–2020 time period. Metadata in the form of statistical tables, maps and figures are provided along with observation time series.

The GRQA dataset, supplementary metadata and figures are available for download on the DataCite- and OpenAIRE-enabled Zenodo repository at <https://doi.org/10.5281/zenodo.5097436> (Virro et al., 2021).

1 Introduction

Human-driven loads of nutrients to aquatic ecosystems have become the main driver of eutrophication in waterways and coastal zones (Desmit et al., 2018; Sinha et al., 2019). Agricultural production is already one of the major forces behind environmental degradation (Foley et al., 2011), and population growth is increasing that pressure (Mueller et al., 2012). The use of nitrogen (N) and phosphorus (P) fertilizers to increase agricultural productivity is predicted to increase 3-fold by 2050 unless more efficient fertilizer use can be im-

plemented (Tilman et al., 2001). At the same time, it has been estimated that “globally, over 3 billion people are at risk of disease because the water quality of their water source is unknown due to a lack of data” (UN-Water, 2021). In order to achieve the UN Sustainable Development Goal (SDG) 6, we need better understanding of our water resources and water quality. Monitoring and modeling the hydrochemical properties of rivers is essential for understanding and mitigating water quality deterioration due to agricultural and industrial non-point-source pollution (Krysanova et al., 1998; Leon et al., 2001; Wu and Chen, 2013). Modeling of different

water quality indicators such as nutrients (Caraco and Cole, 1999; He et al., 2011), carbon compounds (Evans et al., 2005; Hope et al., 1994), sediments (Choubin et al., 2018; Ouyang et al., 2018) and oxygen (Radwan et al., 2003; Singh et al., 2009) gives valuable understanding of hydrochemical cycles and enables us to estimate the effect of human influence on them.

Traditional approaches to water quality modeling consist of applying bottom-up, physically based models on the catchment level (Wellen et al., 2015). Calibration and validation data in the form of water quality observations used when developing the model and verifying its performance are usually gathered through in situ observations and, more recently, automated sensor networks. Although airborne remote-sensing-based data acquisition methods have been successfully used to supplement field data for lakes (Chen and Quan, 2011; Toming et al., 2016), applying those methods is only viable in the case of rivers with a large enough surface area (Olmanson et al., 2013). Therefore, improving the river water quality data's spatial and temporal coverage with remote sensing is limited. Significant progress has been made in improving the technical capabilities and lowering the installation and maintenance costs of the field sensors, but the spatial and temporal coverage of observation sites remains an issue (Pellerin et al., 2016).

In order to improve the spatial coverage of water quality and hydrological data, different solutions have been used in predictive hydrological mapping. Until recently, a common approach for predicting water quality and hydrological phenomena in ungauged catchments has been the application of already-existing process-based models to catchments with similar characteristics (Hrachowitz et al., 2013; Strömqvist et al., 2012; Wood et al., 2011). These physical models usually require extensive calibration along with location-specific knowledge, which limits the wider applicability and spatial upscaling that can be done (Abbaspour et al., 2015; McMillan et al., 2012).

Recently, advances in implementing machine learning (ML) methods in hydrology have given rise to a new, data-driven approach to hydrological modeling (Mount et al., 2016). A comparison of physically based and ML approaches has shown that ML methods can achieve a similar accuracy to the physically based ones and outperform them when describing nonlinear relationships (Chau, 2006; Ouali et al., 2017; Papacharalampous et al., 2019). The recent advance of so-called physics-guided ML, which entails combining process-based models with ML methods, is likely to become more applicable in the near future as well (Kratzert et al., 2019; Shen et al., 2018; Marzadri et al., 2021).

Nevertheless, a major problem related to large-scale predictive hydrological modeling has been the lack of available observation data with a good spatiotemporal coverage (Bierkens, 2015). This has affected the reproducibility of previous studies and the potential improvement of existing models (Blöschl et al., 2019; Meals et al., 2010; Stagge et al.,

2019). In addition to the observation data themselves, insufficient or poor-quality metadata have also discouraged researchers from integrating the already-available datasets. Here, ambiguities in supplementary metadata such as parameter names, units and methods of measurement have limited the use of open data for large-scale water quality modeling purposes (Archfield et al., 2015; Hutton et al., 2016; Sprague et al., 2017). Therefore, improving both the availability and quality of open water quality data would increase the potential to implement predictive modeling on a global scale. Global ML models have been already successfully used for discharge modeling (Beck et al., 2015; Gudmundsson and Seneviratne, 2015), and recent years have seen the publication of global discharge datasets (Do et al., 2018; Harrigan et al., 2020). The publication of global and continental datasets (Hartmann et al., 2014; Read et al., 2017) could make ML methods applicable for large-scale water quality modeling as well (Shen et al., 2020). However, issues related to a lack of training and validation data due to general data scarcity affect model accuracy and, therefore, limit the further adoption of ML for global water quality predictions (Chen et al., 2020).

We aim to address the aforementioned issues by presenting the novel Global River Water Quality Archive (GRQA) by integrating and harmonizing five different global and regional datasets. The resulting dataset has combined observation data for 42 different forms of some of the most important water quality parameters relevant for nutrients (e.g., water temperature and oxygen, phosphorus, nitrogen and carbon compounds). Supplementary metadata and statistics are provided with the observation time series to improve the usability of the dataset. An extensive data catalog with maps showing the spatiotemporal coverage and graphs describing the distribution of all 42 parameters as the supplementary material of the study are presented in the Supplement. We report on developing a harmonized schema and reproducible workflow that can be adapted to integrate and harmonize further data sources. In addition, we provide recommendations for improving the multi-source water quality data compilation, especially focusing on the metadata quality and adhering to the FAIR data principles (Wilkinson et al., 2016). We conclude our study with a call for action to extend this dataset and hope that the provided reproducible method of data integration and metadata provenance shall lead as an example.

2 Data

A total of five data sources were used to compile the GRQA, with two being at the global, one at the regional and two at the national level (Table 1). All datasets with the exception of GEMStat are publicly available to download online as CSV or Excel file packages. GEMStat data can be requested via email. The number of available observation sites was highly dependent on the source, with the Water Quality

Portal (WQP) maintained by the United States Geological Survey (USGS) having the most sites. Files used during the creation of GRQA are listed in Table 2.

2.1 CESI

The first dataset included in GRQA originated from the Canadian Environmental Sustainability Indicators (CESI) program operated by Environment and Climate Change Canada (ECCC), which is a Canadian governmental department responsible for coordinating environmental policies and programs. CESI consists of water quality measurements collected by federal, provincial and territorial monitoring programs from Canadian rivers from the 2002–2018 time period (Environment and Climate Change Canada, 2020). CESI data are mainly focused on heavy metals, so out of the 42 parameters included in GRQA, only 8 were available in CESI (Table 1). It is the smallest of the five source datasets with the site count ranging from 2 to 77 per parameter. Mean time series length per site is approximately 13 years, and the average number of observations per site is 145.

2.2 GEMStat

The Global Freshwater Quality Database GEMStat is hosted by the International Centre for Water Resources and Global Change (ICWRGC) and provides inland water quality data within the framework of the GEMS/Water program of the United Nations Environment Programme (UNEP). GEMStat contains over 7 million samples from approximately 5700 sites in 75 countries. The data were obtained through a custom request to their data portal (United Nations Environment Programme, 2020).

Approximately 500 water quality parameters were available in the GEMStat database, out of which 32 were used when compiling GRQA (Table 1). Observations cover the period 1950–2020, and the mean observation count per parameter is approximately 41. Mean time series length per site is 9 years. The site count per parameter ranges from less than 10 (dissolved and total carbon) to 4274 (total phosphorus).

2.3 GLORICH

The GLObal RIver CHemistry (GLORICH) database (Hartmann et al., 2014) is a collection of hydrochemical data from more than 1.27 million observations and more than 18 000 sampling locations across the globe. The samples originate from various environmental monitoring programs and scientific literature.

Out of 47 water quality parameters available in the raw data, 26 were chosen to be included in the GRQA (Table 1). The samples cover the time period of 1942–2011, but the length of the time series is dependent on the parameter. Mean time series length per site is less than a decade for all parameters. The number of available sites per parameter ranges from

just 4 (particulate organic nitrogen) to 9728 (dissolved inorganic phosphorous). The dataset can be downloaded from Pangaea (Hartmann et al., 2019).

2.4 Waterbase

Waterbase is the generic name given to the European Environment Agency's (EEA) databases on the status and quality of Europe's rivers, lakes, groundwater bodies, and transitional, coastal and marine waters (European Environment Agency, 2020). The database is compiled from data sent by the national European water agencies involved in the Water Framework Directive (WFD).

Over 600 water quality parameters are included in the full dataset, out of which 15 matched those of GRQA (Table 1). Out of all source datasets, Waterbase had the shortest time series with observations covering only the period 2008–2018. The maximum site count per parameter is 1976, while there were on average only around 19 observations per site.

In May 2020, the ICWRGC announced that parts of Waterbase had also been added to the GEMStat database (United Nations Environment Programme, 2020). However, only sites with more than 3 years of data were included in this update. As mean time series length per site was only 1.4 years in Waterbase, a significant number of sites were left out, which is why we decided to include Waterbase separately in GRQA. Although it is likely that there were many observations which appeared both in GEMStat and Waterbase, the duplication detection procedure discussed in Sect. 3.3 should have identified them.

2.5 WQP

USGS, the US Environmental Protection Agency (EPA) and the National Water Quality Monitoring Council developed the Water Quality Portal (WQP), which is so far the largest standardized water quality database (Read et al., 2017; United States Geological Survey, 2020). Although the portal also includes data from a few other countries (e.g., Mexico, Pacific Islands) associated with the National Water Information System (NWIS) network, only a very limited amount of non-US samples were available. For this reason, only US national data were selected to be added to GRQA.

Due to the size of the source dataset, the full set of parameters could not be downloaded at once. Therefore, a scripted download procedure was used to retrieve water quality samples and their corresponding sampling sites separately per parameter. In the case of temperature, the data had to be further divided by state. Unlike other source datasets used in the study, the WQP often had multiple versions of the same parameter available under separate codes in case the parameter had been measured in different units, e.g., using different methods. The final count of parameters used for GRQA was 37 (Table 1).

Table 1. Source datasets used for compiling GRQA with their total number of observations, parameters and timeframe length in GRQA. All datasets were retrieved on 16 November 2020.

Dataset	Name	Data provider	Observations	Timeframe	Parameters (source/GRQA)	Site count range	Mean time series length per site years	Mean observation count per site
			<i>n</i>		<i>n/n</i>	<i>n</i>		<i>n</i>
CESI	Water quality in Canadian rivers	Environment Canada	30 457	2002–2018	8/42	2–77	12.9	145
GEMStat	Global Freshwater Quality Database	International Centre for Water Resources and Global Change	2 094 598	1950–2020	32/42	7–4274	9.2	77
GLORICH	GLObal RIver Chemistry database	Institute of Geology of the University of Hamburg	3 231 797	1942–2011	26/42	4–9728	4.1	41
Waterbase	Waterbase – Water Quality	European Environment Agency	306 332	2008–2018	15/42	4–1976	1.4	19
WQP	USGS Water Quality Portal	Environmental Protection Agency	8 689 335	1898–2020	37/42	1–59 000	3.4	25

Table 2. Source dataset files used for compiling GRQA. WQP sites and observations were downloaded separately for each parameter, and file names were assigned during the process.

File name	Size (MB)	Rows	Description	Sheet name	Source
wqi-federal-raw-data-2020-iqe-donnees-brutes-fed.csv	171.5	314 867	Observation data		CESI
data_request.xls	2.4	5419	Site data	Station_Metadata	GEMStat
data_request.xls	2.4	30	Parameter data	Parameter_Metadata	GEMStat
data_request.xls	2.4	311	Method data	Methods_Metadata	GEMStat
pH.csv	21.9	372 211	Observation data		GEMStat
Carbon.csv	19.2	337 928	Observation data		GEMStat
Nitrogen.csv	65.1	1 052 823	Observation data		GEMStat
Phosphorus.csv	24.3	386 113	Observation data		GEMStat
Oxygen_Demand.csv	20.1	331 617	Observation data		GEMStat
Solids.csv	11.8	201 628	Observation data		GEMStat
Water_Temperature.csv	23.9	370 335	Observation data		GEMStat
Oxygen.csv	30.6	488 749	Observation data		GEMStat
Sampling_Locations_v1.shp	0.4	15 553	Site point data		GLORICH
sampling_locations.csv	1.6	18 897	Site name data		GLORICH
catchment_properties.csv	10.2	15 514	Catchment data		GLORICH
hydrochemistry.csv	273.3	1 274 102	Observation data		GLORICH
Waterbase_v2019_1_S_WISE6_SpatialObject_DerivedData.csv	15.1	62 288	Site data		Waterbase
ObservedProperty.csv	0.2	888	Observation data		Waterbase
Waterbase_v2019_1_T_WISE6_DisaggregatedData.csv	10 019.2	39 121 790	Observation data		Waterbase
WQP_*_sites.csv	2543	9 467 369	Site data		WQP
WQP_*_obs.csv	2749.8	10 088 212	Observation data		WQP

The longest time series of source datasets is present in the WQP, with some dating back to 1898. However, the average time series length per station is just over 3 years. Like GEMStat, WQP is still being updated, so most parameters have their latest observations in 2020. The site count ranges from a single station (dissolved inorganic nitrogen) to 59 000 per parameter (total suspended solids).

3 Methodology

The GRQA compilation workflow was divided into three parts: (1) the preprocessing stage involved converting observation data from the five sources into a common format and harmonizing the corresponding metadata; (2) preprocessed data were merged by parameter, after which outliers and time series characteristics were detected; (3) duplicate detection was conducted in the last processing step. The Pandas (McK-

inney, 2010), GeoPandas (Jordahl et al., 2020) and NumPy (Harris et al., 2020) Python libraries were used throughout all data processing stages.

3.1 Source data preprocessing

3.1.1 Parameter selection

The parameters included in GRQA cover the four groups of water quality indicators outlined in the introduction: nutrients, carbon, sediments and oxygen (Table 7). GLORICH was used as a reference for parameter selection due to its being one of the two global source datasets and having the least amount of discrepancies within source data, i.e., each GLORICH parameter had a single matching code, unit, etc.

3.1.2 Parameter harmonization

Preliminary analysis showed that there were ambiguities in the parameter names, codes, units and chemical forms in the different source datasets, which has been identified as a recurring issue when dealing with multi-source water quality data (McMillan et al., 2012; Sprague et al., 2017). For this reason, lookup tables were created for each of the source datasets (**_code_map.csv*) to use as guides in the following processing stages (Table 3). The purpose of the schemas was to match parameter codes and other metadata with the versions used later in the GRQA. For most parameters, this could be done based on the literal names, remarks and descriptions in the metadata. Relevant literature and online resources were consulted for more ambiguous scenarios. One such example was total suspended solids (TSSs), which can also be reported as suspended particulate matter (SPM) (Neukermans et al., 2012). When a reliable decision could not be made (e.g., biological oxygen demand as BOD vs. BOD5), the parameters were kept separate.

3.1.3 Unit conversion

Units of measurement were harmonized along with other metadata. All parameters except temperature (°C), pH and dissolved oxygen (%) were converted into milligrams per liter (mg L^{-1}), which was the most prevalent unit in source data. When units were converted, observation values had to be changed as well. This was done by calculating conversion constants which were based on both the magnitude of the source unit (e.g., $\mu\text{g L}^{-1}$ vs. mg L^{-1}) and the reported chemical form of the parameter. The latter affected nitrite (NO_2), nitrate (NO_3) and ammonium (NH_4) the most as these parameters had a variety of forms in the source data that were all converted into corresponding nitrogen versions ($\text{NO}_2\text{-N}$, $\text{NO}_3\text{-N}$ and $\text{NH}_4\text{-N}$). In some cases, the chemical form could be identified from the source unit (e.g., mg NL^{-1} or $\text{mg NO}_3\text{ L}^{-1}$), while others were detected by examining parameter names and method descriptions (e.g., “Nitrate, reported as nitrogen”). Where possible, additional informa-

tion about these missing forms was collected from proxy sources, such as other similar datasets (e.g., Börker et al., 2020 in the case of GLORICH). These references have been included in the *form_ref* column in corresponding lookup tables (**_code_map.csv*). For other nitrogen (TKN, TN, etc.), all carbon (DOC, TC, etc.) and phosphorus (TP, TIP, etc.) parameters (see Table 7 for abbreviations), the chemicals were assumed to be either N, C or P even if not reported because there is only one common element in the molecule (Sprague et al., 2017). GLORICH was the only source dataset which also needed conversion constants for carbon and phosphorus parameters as they had been reported as micromoles per liter ($\mu\text{mol L}^{-1}$). All WQP units matched those intended to be used for GRQA, so no conversion was needed. The formula for conversion constants was

$$x_2 = \frac{x_1 \times M_{x_2}}{n \times M_{x_1}}, \quad (1)$$

where x_1 and x_2 are observation values before and after conversion, M is the corresponding molar mass, and n is the magnitude difference between source and converted unit. Some examples of unit conversion are given in Table 4. The full list of all unit conversion procedures is given in Appendix A (Table A1).

3.1.4 Site ID duplication

There were some instances of duplicated site IDs in GLORICH (2 site pairs) and Waterbase (101 pairs) source data, which meant that joining observations with sites would have created duplicate time series as well. Site ID duplicates could indicate that there had been small shifts in the site location or that the site had been closed and reinstated at some point. If the distance between the duplicate pairs was less than a kilometer, only the first instance was retained in the output table. When distance was greater than a kilometer, both instances were removed since metadata that could be used to make a decision (e.g., when the site first opened) were not available. Finally, all duplicate pairs were exported as separate files (e.g., *GLORICH_dup_sites*).

3.1.5 Coordinate conversion

CESI and WQP originally had the site coordinates in the North American Datum of 1983 (NAD83). The pyproj (Snow et al., 2020) Python library was used for converting the North American site coordinates into World Geodetic System 1984 (WGS84), which was the coordinate system chosen for the GRQA.

3.1.6 Observation data filtering

Preliminary cleaning included the removal of observations of negative, missing or low-quality values. In this case, low quality refers to measurements that were flagged as either

Table 3. Summary table of lookup table attributes.

Attribute name	Description	Data type
source_param_code	Parameter code in source dataset	String
source_param_code_meta	Additional code specification used for CESI	String
param_code	Parameter code in GRQA	String
source_param_name	Parameter name in source dataset	String
param_name	Parameter name in GRQA	String
source_param_form	Parameter chemical form in source dataset	String
param_form	Parameter chemical form in GRQA	String
form_ref	Parameter form reference	String
source_unit	Parameter unit in source dataset	String
divisor	Divisor applied to the observation value	Float
multiplier	Multiplier applied to the observation value	Float
conversion_constant	Unit conversion constant calculated based on divisor and multiplier and applied to the observation value	Float
unit	Parameter unit in GRQA	String
source	Source dataset name	String

Table 4. Examples of unit conversion from the chemical form in source data to the GRQA version; x_1 and x_2 are observation values before and after conversion, respectively.

Parameter code	Source	Form	Source form	Unit	Source unit	x_1	M_{x_2}	n	M_{x_1}	x_2
TAN	CESI	N	NH ₃	mg L ⁻¹	mg L ⁻¹	0.106	14.007	1	17.031	0.087
NO2N	GEMStat	N	NO ₂	mg L ⁻¹	mg L ⁻¹ NO ₂	0.024	14.007	1	46.005	0.007
NO3N	GLORICH	N	NO ₃	mg L ⁻¹	μmol L ⁻¹	210.268	14.007	1000	62.004	0.048
NH4N	Waterbase	N	NH ₄	mg L ⁻¹	mg L ⁻¹	0.063	14.007	1	18.039	0.049

coming from unreliable sources or having any kind of literal quality assessment flag in the source data (e.g., “poor quality”). Additionally, observations marked as below (<) or above (>) detection limit in source data were flagged as such in GRQA as well (column *detection_limit_flag*). Observations originating from unreliable sources or being otherwise suspect (e.g., unvalidated) were omitted. Three source datasets (GEMStat, GLORICH and Waterbase) had this type of a quality evaluation included in the metadata. Observations from sites marked as “not for publication” due to national legislation in Waterbase were also not included in GRQA.

3.1.7 Filtration information

When possible, supplementary information about whether a sample was filtered or unfiltered was retained as filtration can affect the sample values (Sprague et al., 2017). This information was usually available in a separate metadata column. Both “filtered” and “dissolved” were used depending on the source. GRQA includes the dissolved versions of certain parameters (total nitrogen, total phosphorus and Kjeldahl nitrogen) which originally did not exist as separate parameters in Waterbase and WQP. In those cases, the filtered and dissolved observations of TN, TP and TKN in the two

datasets were treated as the corresponding dissolved forms (TDN, TDP, DKN) in GRQA.

3.1.8 Time and date processing

Observations could have invalid timestamps due to formatting or entry errors, so a validity check was included in the preprocessing scripts. Dates were tested against the presumed source format, and observations with incorrectly formatted or implausible dates were removed. The source datasets used different date formats, which were all converted into a common one (%Y-%m-%d). When possible, observation time was extracted as well. A default value (00:00:00) was used to fill missing information. Time zone information was only able to be extracted from the WQP. Other sources lacked time zone information, so it was not possible to determine whether the recorded timestamp was in local or coordinated universal time (UTC), and the time given is up to the user to interpret.

3.1.9 Other metadata

If available, metadata about the upstream basin area, its unit and the name of the greater drainage region of the site were included in GRQA. Additional information about methods used or other available observation remarks in the source data was also retained. The metadata depended on the source and

were available only sporadically and could not be concatenated in a reasonable way between the datasets, so the information is given in the GRQA for each source separately in the format of *source_meta_sourcecolumnname* (e.g., *GEMSTAT_meta_Analysis Method Code*). Here, the source column names were kept as they appear in raw data, e.g., spaces were not replaced with underscores.

3.2 Outlier treatment, time series availability and continuity

3.2.1 Time series availability and continuity

The analysis of the statistics generated during preprocessing showed that most of the time series extracted from the source datasets are very discontinuous. For example, the mean time series length per site for total phosphorus (TP) in GEMStat was 6.6 and in GLORICH 4.9 years, while the mean observation count per site was only 57.7 and 52.4, respectively. This means that many sites have observations at a monthly time step at best. Similar findings have been previously reported for WQP time series (Read et al., 2017; Shen et al., 2020).

In order to illustrate the suspected temporal fragmentation in observation data, monthly availability and monthly continuity statistics appropriated from the strategy used by Crochemore et al. (2019) were calculated for each site in each of the merged parameter time series. Both characteristics can give insight into the granularity of the time series and can affect the applicability of different modeling methods. Monthly availability of observation data was defined as the ratio between number of months with at least one observation and the total number of months a particular site had any observations. A ratio of 1.0 would mean that there was at least one observation in every month of the time series. Monthly continuity was calculated as the ratio between the longest period of consecutive months with any measurements and the length of time series in months. Here, a ratio of 1.0 would mean that there were no months without observations, and the time series is continuous on a monthly level. The resulting characteristics were added as columns in the output files.

3.2.2 Outlier flagging

Water quality modeling often involves dealing with numerous outliers and uncertainties in observation data, particularly when integrating time series from multiple sources (McMillan et al., 2012; Sprague et al., 2017). Due to the differences in environmental conditions and water regimes, the potential range of observation values can vary a lot between catchments. Although extreme outliers caused by faulty equipment or data entry errors can sometimes be detectable by examining distribution plots, it is often difficult to decide whether an outlier is an error or not. For example, sudden spikes in observation time series can be caused by events such as accidental fertilizer spills in the waterway or a

cow getting entrapped in an in-stream wetland (Hughes et al., 2016), which can have short-term effects on water quality and, therefore, should not be removed from data. However, flagging outliers can still help researchers troubleshoot potential issues at the modeling stage.

For this reason, no observations were omitted from the time series, and two flags associated with outliers were added to the output tables instead. First flag (*obs_iqr_outlier*) shows whether an observation was deemed to be an outlier by the interquartile range (IQR) test. IQR is defined as the difference between the third (Q3) and first (Q1) quartile. All values greater than $Q3 + 1.5 \times IQR$ or less than $Q1 - 1.5 \times IQR$ are considered outliers. The second flag (*obs_percentile*) was an indicator (0.0–1.0) showing which percentile a particular observation belongs to. Histograms along with box and whisker plots were used to visually show the range and distribution of the parameter observations. The plots were produced for every parameter and are included in the GRQA data repository.

3.3 Duplicate observation detection

The global datasets (GEMStat and GLORICH) used in this study had at least partial spatial overlap with the other three sources, which means that merging could have created duplicate sites in the GRQA. Contrary to site ID duplicates within the same dataset discussed in Sect. 3.1, site duplicates from different sources would likely also have different IDs. Therefore, rather than comparing ID information, the duplicates had to be identified by spatial proximity and time series similarity. Similar to procedures described in Sect. 3.2, duplicate detection was done separately for each parameter.

The first stage of duplicate detection was clustering sites based on their geographic location. The DBSCAN (density-based spatial clustering of applications with noise) algorithm (Xu et al., 1998) from the scikit-learn Python library (Pedregosa et al., 2011) was used to create clusters of sites within a 1 km radius of each other, which is the approximate accuracy of around two decimal points in latitude and longitude degrees. There does not seem to be a consensus for assigning this search radius for duplicate detection, and the assessment of spatial proximity depends on the subjective threshold set by authors. For example, the Global Streamflow Indices and Metadata Archive (GSIM) streamflow dataset (Do et al., 2018) used a radius of 5 km for selecting potential duplicate gauging stations. The 1 km radius was chosen to avoid having too many false positives (e.g., in the case of small headwater catchments) to evaluate in the second stage of deduplication (RMSE calculation). A major advantage of DBSCAN compared to similar density-based clustering methods is that the algorithm can be run without determining a priori the number of output clusters (Birant and Kut, 2007). In addition, DBSCAN has been shown to be more applicable than others when dealing with large-scale datasets (Khan et al., 2014; Parimala et al., 2011).

Although there are time series similarity detection methods that can be applied to irregular time series and can handle some degree of discontinuity, the focus of those methods is on misalignment of the time of observations rather than differences in the pattern of time series gaps (Berndt and Clifford, 1994). Therefore, it is likely that GRQA time series are too fragmented for these advanced methods to yield reliable results. A conservative approach based on root-mean-square error (RMSE) was chosen here instead. Output site clusters were converted into unique site pairs so that all sites within a cluster could be compared to one another (e.g., a cluster of four would yield six unique ID pairs). Site ID pairs were then used to extract corresponding time series from observation data. Only observations made on matching dates were used for calculating the RMSE, and only pairs in which RMSE was equal to zero were considered as potential duplicates. Finally, the duplicates were exported into separate CSV files (e.g., *TP_dup_obs.csv*) along with relevant metadata to help the user decide whether the sites can be considered duplicates (Table 5). A high number of matching dates with the same observation value (column *date_match_count*) would indicate a higher likelihood of duplication.

4 Results

4.1 GRQA data model and descriptive overview

The GRQA dataset consists of observation time series for 42 different water quality parameters provided in tabular form as CSV files. Each of the observation files is accompanied by corresponding metadata files (tables and images) describing the spatial and temporal characteristics of the time series.

GRQA is made up of the following files (Fig. 1):

- a data catalog (*GRQA_data_catalog.pdf*) with maps showing the spatiotemporal coverage and graphs describing the distribution of all 42 parameters along with a README file describing the dataset structure
- water quality observation time series files (named *paramcode_GRQA.csv*)
- GRQA metadata (folder *meta*) containing descriptive statistics (*GRQA_param_stats.csv*) and duplicate observation files (*source_dup_obs.csv*) when relevant
- the set of overview figures (folder *figures*) containing
 - histograms (*paramcode_GRQA_hist.png*) and box plots (*paramcode_GRQA_box.png*) showing the distribution of observation values by source dataset
 - maps showing the spatial distribution of the observations by source (*paramcode_GRQA_spatial_dist.png*)
 - maps showing the median observation values of sites (*paramcode_GRQA_median.png*)

maps showing the monthly availability (*paramcode_GRQA_availability.png*) and continuity (*paramcode_GRQA_continuity.png*) of the observations.

The five source datasets are also included in the GRQA data package. Folder *source_data* includes

- the *raw* folder with downloaded source files and harmonization schemas used in the preprocessing stage (*source_code_map.csv*) for each source dataset along with the original units (*source_units.csv*)
- the *sourcename/processed* folder containing summary statistics of observation values by parameter for each source dataset before (*paramcode_source_raw_stats.csv*) and after (*paramcode_source_processed_stats.csv*) processing along with information about the number of missing values (*source_missing_values.csv*) and source file size (*source_file_info.csv*)
- when relevant, *processed/meta* also including duplicate site ID files (*source_dup_sites.csv*).

The structure of GRQA observation files is given in Table 6. In addition to the attributes outlined in Sect. 3, the extracted metadata also include information about the upstream basin and drainage region of the observation site. It has to be noted that the availability of this information was dependent on both the source (i.e., not present in CESI and Waterbase) and the observation site itself and is therefore available only sporadically in GRQA as well (Table 6). Parameter codes, names, forms and observation values in GRQA are given as they appeared in source data alongside their harmonized and processed GRQA versions so that end users could assess the validity of conversion and make corrections if needed.

The statistical overview of the parameters included in GRQA is shown in Table 7. The number of sites per parameter ranges from only 7 (DC) up to 68 592 (TSSs). Parameters having more sites generally also have more observations. Parameters with a small number of sites and observations were usually present in only one or two source datasets. For example, dissolved organic phosphorus (DOP) only existed in WQP. Different versions of biochemical and chemical oxygen demand that could not be harmonized based on source metadata were kept separate, although the median value for BOD and BOD5 ended up being equal.

The spatial distribution of water quality observation sites depended on the parameter and is illustrated in Fig. 2 using dissolved oxygen (DO), dissolved organic carbon (DOC), TP and TSSs. These parameters were the largest in terms of number of sites and observations in their corresponding groups (oxygen, carbon, nutrients and sediments). They are also used in the following figures. Some observations that could be made when examining site maps were the following:

Table 5. Summary table of duplicate observation file attributes.

Attribute name	Description	Data type
obs_id_1	Observation ID of first site	String
lat_wgs84_1	Latitude of first site	Float
lon_wgs84_1	Longitude of first site	Float
site_id_1	First site ID	String
site_name_1	First site name	String
obs_value_1	First site observation value	Float
source_1	First site source	String
site_ts_availability_1	First site availability	Float
site_ts_continuity_1	First site continuity	Float
obs_date	Observation date	String
obs_id_2	Observation ID of second site	String
lat_wgs84_2	Latitude of second site	Float
lon_wgs84_2	Longitude of second site	Float
site_id_2	Second site ID	String
site_name_2	Second site name	String
obs_value_2	Second site observation value	Float
source_2	Second site source	String
site_ts_availability_2	Second site availability	Float
site_ts_continuity_2	Second site continuity	Float
date_match_count	Number of matching dates with the same observation value	Int
param_code	Parameter code	String

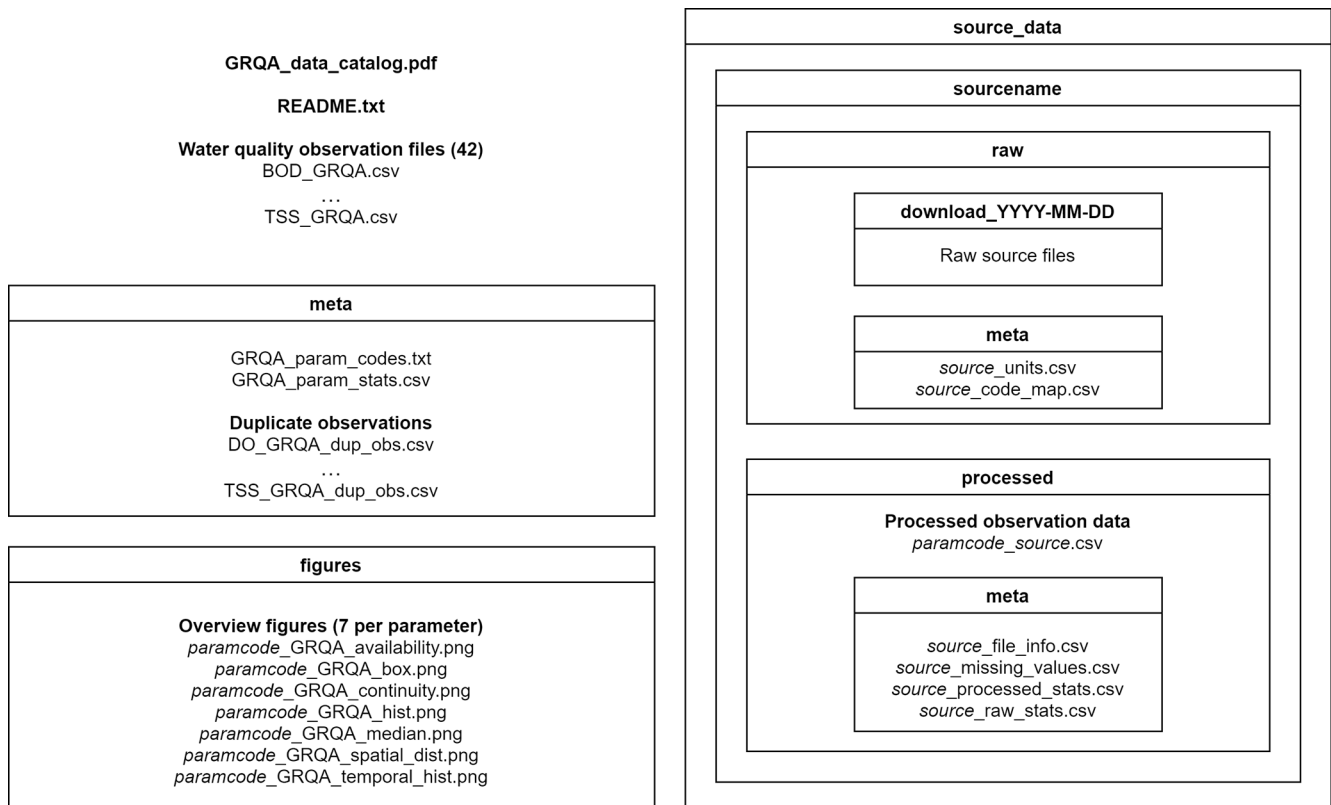


Figure 1. Diagram showing the folder structure and contents of the GRQA dataset.

Table 6. Summary table of output water quality observation file attributes.

Attribute name	Description	Data type
obs_id	Unique observation ID generated by hashing	String
lat_wgs84	Observation site latitude in WGS84	Float
lon_wgs84	Observation site longitude in WGS84	Float
obs_date	Observation date in the %Y-%m-%d format	String
obs_time	Observation time in the %H:%M:%S format	String
obs_time_zone	Observation time zone code	String
site_id	Observation site ID	String
site_name	Observation site name	String
site_country	Observation site country	String
upstream_basin_area	Site upstream basin area	String
upstream_basin_area_unit	Site upstream basin area unit	String
drainage_region_name	Drainage region where site is located in	String
param_code	Parameter code in GRQA	String
source_param_code	Parameter code in source dataset	String
param_name	Parameter name in GRQA	String
source_param_name	Parameter name in source dataset	String
obs_value	Observation value in GRQA	Float
source_obs_value	Observation value in source dataset	Float
detection_limit_flag	Whether a value was flagged as below (<) or above (>) detection limit in source data	String
param_form	Parameter chemical form in GRQA	String
source_param_form	Parameter chemical form in source dataset	String
unit	Parameter unit in GRQA	String
source_unit	Parameter unit in source dataset	String
filtration	Sample filtration information	String
source	Source dataset name	String
obs_percentile	Percentile of the observation value	Float
obs_iqr_outlier	Flag to mark whether observation value is an outlier according to the interquartile range test	String
site_ts_availability	Monthly availability of the time series per site	Float
site_ts_continuity	Monthly continuity of the time series per site	Float
meta	Other observation metadata with a reference to the corresponding source column (e.g., GEMSTAT_meta_Method Description)	String
...	...	

Table 7. GRQA water quality parameter statistics.

Parameter code	Parameter name	Sites	Observations	Median value	Unit	Start year	End year	Outlier %
BOD	Biochemical oxygen demand	2945	163 531	2.627	mg L ⁻¹	1974	2019	13.4
BOD5	Biochemical oxygen demand (BOD5)	13 283	278 629	5.875	mg L ⁻¹	1905	2020	8.3
BOD7	Biochemical oxygen demand (BOD7)	386	5282	2.200	mg L ⁻¹	2013	2018	5.9
COD	Chemical oxygen demand	2769	126 372	22.362	mg L ⁻¹	1974	2019	10.8
CODCr	Chemical oxygen demand (Cr)	671	7350	24.900	mg L ⁻¹	2013	2018	3.4
CODMn	Chemical oxygen demand (Mn)	287	2310	4.600	mg L ⁻¹	2013	2018	2.3
DC	Total dissolved carbon	7	9	4.800	mg L ⁻¹	2000	2001	0
DIC	Dissolved inorganic carbon	969	30 633	12.266	mg L ⁻¹	1968	2020	3.5
DIN	Dissolved inorganic nitrogen	119	7822	4.200	mg L ⁻¹	1998	2019	2.6
DIP	Dissolved inorganic phosphorus	9931	612 922	0.040	mg L ⁻¹	1942	2017	13.3
DKN	Dissolved Kjeldahl nitrogen	2820	80732	0.347	mg L ⁻¹	1973	2020	6.5
DO	Dissolved oxygen	48 072	1 487 724	8.835	mg L ⁻¹	1898	2020	2.2
DOC	Dissolved organic carbon	14 799	413 328	2.804	mg L ⁻¹	1968	2020	6.8
DON	Dissolved organic nitrogen	10 811	163 630	0.371	mg L ⁻¹	1951	2020	8.1
DOP	Dissolved organic phosphorus	142	899	0.010	mg L ⁻¹	1971	2003	8.7
DOSAT	Dissolved oxygen saturation	34 949	953 274	92.164	%	1898	2020	8.7
NH4N	Ammonium nitrogen	11 372	651 850	0.027	mg L ⁻¹	1942	2018	15.1
NO2N	Nitrite nitrogen	30 902	720 944	0.010	mg L ⁻¹	1900	2020	12.7
NO3N	Nitrate nitrogen	45 422	1 229 584	0.468	mg L ⁻¹	1900	2020	11.1
PC	Particulate carbon	2898	51 049	0.908	mg L ⁻¹	1995	2020	11
pH	pH	27 577	1 372 794	6.886	pH	1900	2020	14.1
PIC	Particulate inorganic carbon	1095	9196	0.060	mg L ⁻¹	1974	2020	14
PN	Particulate nitrogen	2996	56 125	0.129	mg L ⁻¹	1981	2020	9.5
POC	Particulate organic carbon	22 910	615 941	1.617	mg L ⁻¹	1900	2020	9.7
PON	Particulate organic nitrogen	28	1111	0.120	mg L ⁻¹	1989	2019	14
POP	Particulate organic phosphorus	12	13	0.020	mg L ⁻¹	1999	2000	7.7
TAN	Total ammonia nitrogen	27 980	717 776	0.065	mg L ⁻¹	1900	2020	13.3
TC	Total carbon	1181	12 338	27.000	mg L ⁻¹	1968	2007	3.3
TDN	Total dissolved nitrogen	968	62 980	0.310	mg L ⁻¹	1972	2020	11.2
TDP	Total dissolved phosphorus	3325	169 297	0.031	mg L ⁻¹	1965	2020	11.3
TEMP	Water temperature	26 860	1 113 471	18.968	°C	1912	2020	9.3
TIC	Total inorganic carbon	1984	23 024	11.833	mg L ⁻¹	1968	2019	3.8
TIN	Total inorganic nitrogen	78	12 951	3.649	mg L ⁻¹	1992	2020	0.8
TIP	Total inorganic phosphorus	1328	42 495	0.026	mg L ⁻¹	1971	2018	13.8
TKN	Total Kjeldahl nitrogen	9418	425 595	0.680	mg L ⁻¹	1962	2020	8.1
TN	Total nitrogen	18 507	575 887	1.329	mg L ⁻¹	1958	2020	11.9
TOC	Total organic carbon	18 032	420 029	4.526	mg L ⁻¹	1958	2020	7.2
TON	Total organic nitrogen	22 799	592 654	0.622	mg L ⁻¹	1900	2020	8.6
TOP	Total organic phosphorus	294	1811	0.030	mg L ⁻¹	1971	2020	11.9
TP	Total phosphorus	44 990	1 914 538	0.105	mg L ⁻¹	1900	2020	11.8
TPP	Total particulate phosphorus	77	5836	0.021	mg L ⁻¹	1978	2019	10.5
TSSs	Total suspended solids	68 592	1 958 429	9.785	mg L ⁻¹	1898	2020	20.5

– Europe and North America are the best represented in the case of all parameters.

– Coverage is also good in Australia, New Zealand, parts of East Asia and Brazil in the case of some of the key parameters (e.g., TP, TN).

– The rest of the world (Africa, most of Asia) only has sporadic coverage.

The temporal distribution of the four parameters is given in Fig. 3. Similar to the spatial distribution, the temporal coverage of observations depended on both source data and parameter, with WQP having the longest and Waterbase the shortest time series. Most of the data from GEMStat are from the

past decade, while GLORICH has a more even observation distribution throughout the time series.

4.2 Statistical characteristics of GRQA observation time series

As mentioned in the previous section, each of the observation files was accompanied by a set of images and tables, giving insight into the characteristics of the observation time series. The structure of tabular summary statistics is shown in Table 8. These files contain some basic statistics (e.g., standard deviation) about observation values per parameter and source. In addition, information about the temporal characteristics of time series (e.g., mean length per site) is given; this can also be important when assessing the suitability of the data for modeling purposes.

The applicability of water quality modeling is greatly affected by the distribution of observation values as a majority of modeling methods require a near-normal distribution. The skewness caused by extreme outliers is a common problem in hydrological modeling. The observations often follow a log-normal distribution, which means that the data often need to be transformed and normalized in order to be usable (Helsel, 1987; Hirsch et al., 1982; Parmar and Bhardwaj, 2014). Similar behavior was also examined in GRQA, in which values of most parameters showed a strong positive skew. This can be seen in histograms (Fig. 4) and box plots (Fig. A1). For illustrative purposes, values determined as outliers by the IQR test have been omitted from the figures. In the case of parameters such as TP and TSSs, the skewness remains even after outlier omission. This is confirmed by the box plots, in which the total range of the values greatly exceeds the median.

Availability (Fig. 5) and continuity (Fig. 6) plots were used to examine the temporal fragmentation of the time series. In general, observations from national sources (CESI and WQP) exhibited slightly higher availability and continuity than others, likely caused by more consistent data acquisition frameworks. No clear spatial pattern emerged from the analysis, meaning that differences in both indicators exist at the site level even within the same country. Due to how the metrics were calculated, shorter time series outperformed longer ones. An example of this is TP in Brazil, where the examined high continuity correlated with very short mean time series length (less than a year). Parameters with very fragmented time series (e.g., TSSs) had only a limited number of sites where observations had been collected consistently throughout the whole time frame.

The GRQA also includes plots of median observation values, which were calculated over the whole time series for each site. Seasonal fluctuations cannot be identified on this aggregation level, so the maps are meant to be only indicative. An example of median plots can be seen in the Appendix (Fig. A2).

5 Discussion

5.1 Limitations and considerations regarding the use of GRQA

Taking into account the aforementioned issues encountered during the compilation of GRQA, certain limitations and potential remaining errors have to be considered when using the dataset for water quality modeling.

5.1.1 Potential errors in unit conversion

As described in Sect. 3, several assumptions had to be made when creating harmonization schemas about the chemical form of certain nitrogen parameters (NO_2 , NO_3 and NH_4). However, if the assumption made based on this limited ancillary information was incorrect, then using the conversion would have been affected as well. For this reason, the source observation values along with source units were retained, and the users can retrace the conversion steps using the harmonization schemas.

5.1.2 Skewness of observation values

The outlier treatment strategy used for GRQA involved only flagging the values based on the IQR test, which means that the skewness illustrated in Sect. 4 still remains. Although the described strong positive skew existed also in source data, potential unit conversion errors could have exaggerated it. As shown by histograms, omitting flagged outliers is not enough to eliminate the skewness in some cases (TP and TSSs), so additional processing could be needed to transform the data into a normal shape. Power transformation methods like the Box–Cox transformation (Box and Cox, 1964) could be used to further minimize skewness. It is likely that some of the most extreme outliers are caused by data entry errors or equipment malfunction rather than events such as agricultural spills. For setting thresholds to determine whether a value is illogical or not, more sophisticated outlier detection methods based on some general freshwater quality guidelines (Enderlein et al., 1996) could perhaps be used to further filter the observation values.

5.2 Suggestions for improving multi-source water quality data compilation

5.2.1 Metadata quality

When merging datasets from different sources, most of the complications stemmed from inadequate metadata of water quality observations, such as ambiguous parameter names and codes, and missing details on the chemical forms of parameters. This information would be integral for harmonizing units and observation values. The terms used for indicating the filtration status of samples are often dependent on the interpretation of the authors (total vs. unfiltered, dissolved vs.

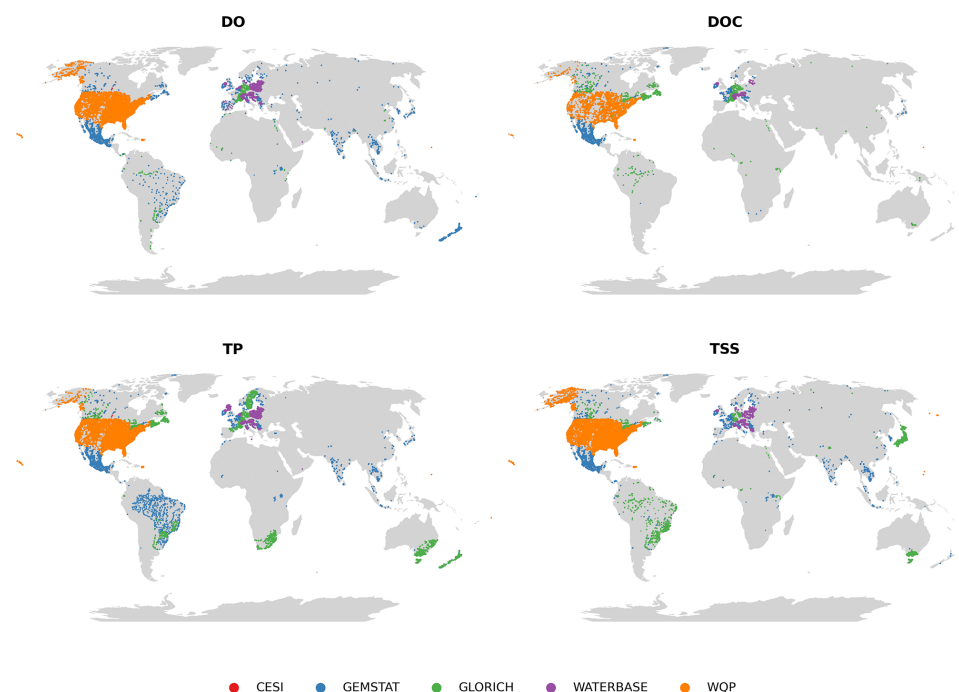


Figure 2. Distribution of observation sites for dissolved oxygen (DO), dissolved organic carbon (DOC), total phosphorus (TP) and total suspended solids (TSSs).

Table 8. Summary table of observation time series statistics file attributes.

Attribute name	Description	Data type
source_param_code	Parameter code in source dataset	String
param_code	Parameter code in GRQA	String
param_name	Parameter name in source dataset	String
source_param_form	Parameter form in source dataset	String
param_form	Parameter form in GRQA	String
source_unit	Parameter unit in source dataset	String
unit	Parameter unit in GRQA	String
count	Total number of observations	Int
min	Minimum observation value	Float
max	Maximum observation value	Float
mean	Mean observation value	Float
median	Median observation value	Float
SD	Standard deviation of observation values	Float
min_year	Time series start	Int
max_year	Time series end	Int
ts_length	Total time series length per parameter	Float
site_count	Total number of sites per parameter	Int
mean_obs_count_per_site	Mean observation count per site	Float
mean_ts_length_per_site	Mean time series length in years per site	Float

filtered), which can affect results when merging (McMillan et al., 2012; Sprague et al., 2017). The annotation of suspect or incomplete data is another aspect of good-quality meta-data (Gudivada et al., 2017). Internal quality control measures such as the ones in GEMStat and Waterbase would help

the end user in the data cleaning stage and eliminate some of the outliers.

The following aspects should be considered to make multi-source data harmonization more efficient in the future:

- Parameter forms should be reported with the units.

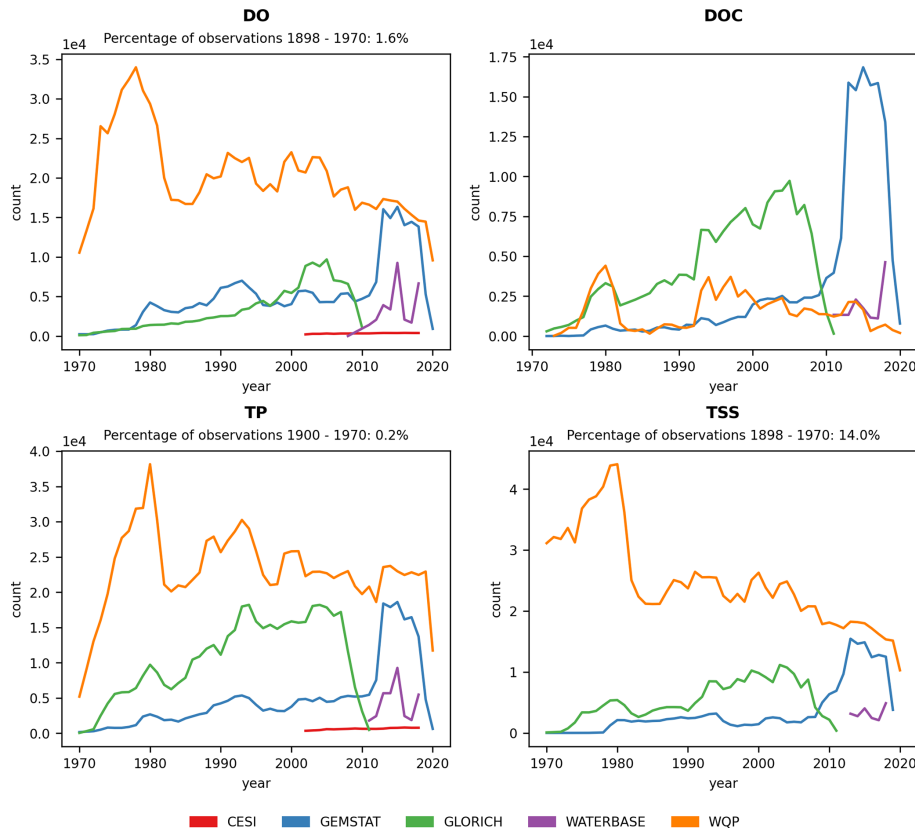


Figure 3. Temporal distribution of observations for dissolved oxygen (DO), dissolved organic carbon (DOC), total phosphorus (TP) and total suspended solids (TSSs) for the period 1970–2020. The percentage of observations before the period shown on the plot is given for each parameter. Only seven observations (1.69×10^{-5} %) existed for DOC in the 1968–1970 period.

- The filtration status of the samples should be reported, and the terms filtered and unfiltered should be preferred as opposed to the more ambiguous dissolved and total.
- Machine-readable quality flags as found in GEM-Stat (columns *Value Flags* and *Data Quality*) or Waterbase (columns *resultObservationStatus*, *meta-data_statusCode* and *metadata_observationStatus*) should be added.
- Whether observations are daily or monthly at the source level should be clearly defined.
- Area units (m^2 , km^2 , etc.) should be included when the upstream catchment area of the site is reported.
- Other information about potential errors in the data should be included (potential duplicates, typographical errors, etc.).
- When certain assumptions or decisions are made when harmonizing data from different sources, they should be reported when the data are published.

5.2.2 Spatial and temporal discontinuity

Although spatial coverage of water quality observations in GRQA exceeds that of the existing global datasets (GEMStat and GLORICH), large areas of Africa and Asia are empty. A major reason might be a lack of knowledge and funding to update and extend site networks, particularly in hard-to-reach areas. In addition, not all governments adhere to an open-data policy. Therefore, improving the spatial coverage of water quality data still relies mostly on implementing additional measures to encourage countries to share it in accordance with open-data principles.

The availability and continuity analysis showed that the GRQA time series are fragmented and that significant gaps remain in the data which will negatively affect large-scale modeling performance. These gaps could be caused by both issues with sensor maintenance or technical limitations under certain conditions (e.g., weather) and inconsistencies in the data acquisition practices on the local level. Recently, ML-based solutions for time series augmentation have been used to fill in gaps in historical monitoring data (Gao et al., 2018; Ren et al., 2019). However, this kind of gap filling still requires enough good-quality training data in the existing time series fragments to be effective and can potentially only be

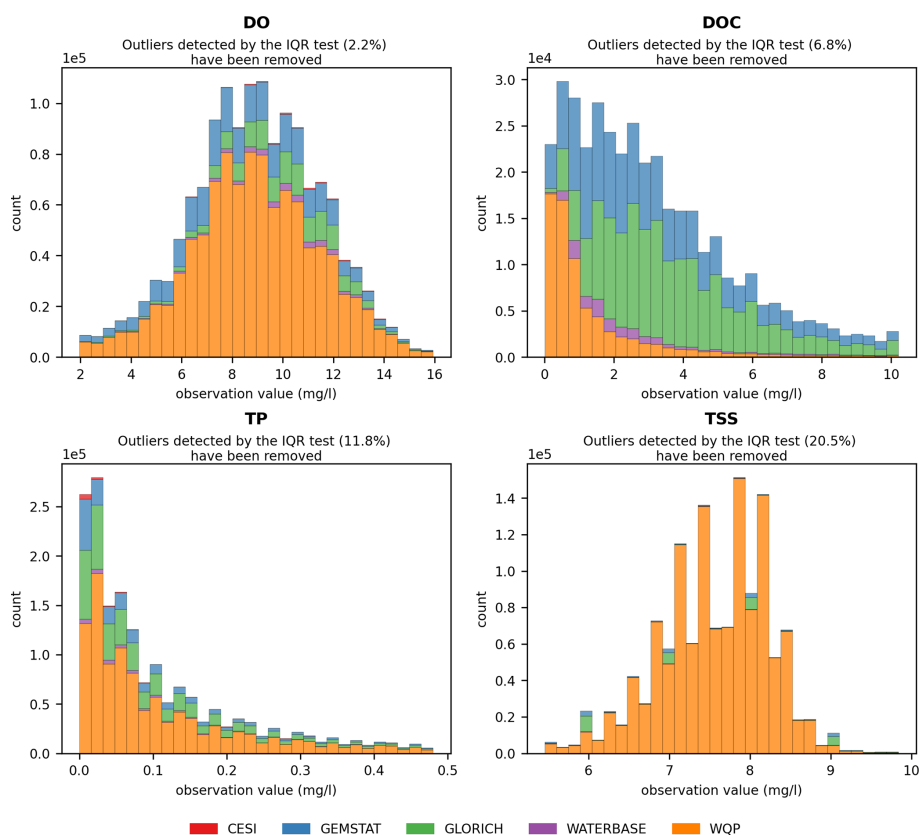


Figure 4. Distribution of observation values for dissolved oxygen (DO), dissolved organic carbon (DOC), total phosphorus (TP) and total suspended solids (TSSs). Outliers determined by the IQR test (Table 7) are not shown on the plot.

of help when improving the temporal rather than spatial coverage.

Another option for improving continuity is using data from one time series to fill in gaps in another. For example, turbidity has been successfully translated into TP and TSS content (Castrillo and García, 2020; Jones et al., 2011). As turbidity data can be acquired at a higher frequency than TP and TSSs, the use of such surrogate parameters can be helpful in data-scarce regions for certain parameters.

5.2.3 General remarks

An important part in improving the spatiotemporal coverage of water quality is raising awareness about the existing datasets (e.g., GEMStat), so that new institutions could join the contributor network and submit their own site data. The continued growth of international collaboration will be vital in improving open global water quality data (Blöschl et al., 2019; Tang et al., 2019). Most of the data collected locally are intended only for regional or national use. Thus, the data are not compatible with those from other countries due to lack of common metadata management practices, with the problems discussed above being a major bottleneck (Hutton et al., 2016; Sprague et al., 2017; Stagge et al., 2019).

Providing those institutions with an example workflow when designing water quality data pipelines, such as the schema recently proposed by Plana et al. (2019), would help them develop their own data management strategy. The workflow used to compile GRQA along with the issues raised in this study will hopefully also help to draw attention to this topic.

6 Code and data availability

The GRQA dataset, supplementary metadata and figures are available for download on the DataCite- and OpenAIRE-enabled Zenodo repository at <https://doi.org/10.5281/zenodo.5097436> (Virro et al., 2021).

The data processing scripts used for the compilation of GRQA are available on Zenodo at <https://doi.org/10.5281/zenodo.5082147> (Virro and Kmoč, 2021).

7 Conclusions

The GRQA dataset was created with the intention to improve the spatiotemporal coverage of previously available open water quality data and provide an example workflow for multi-

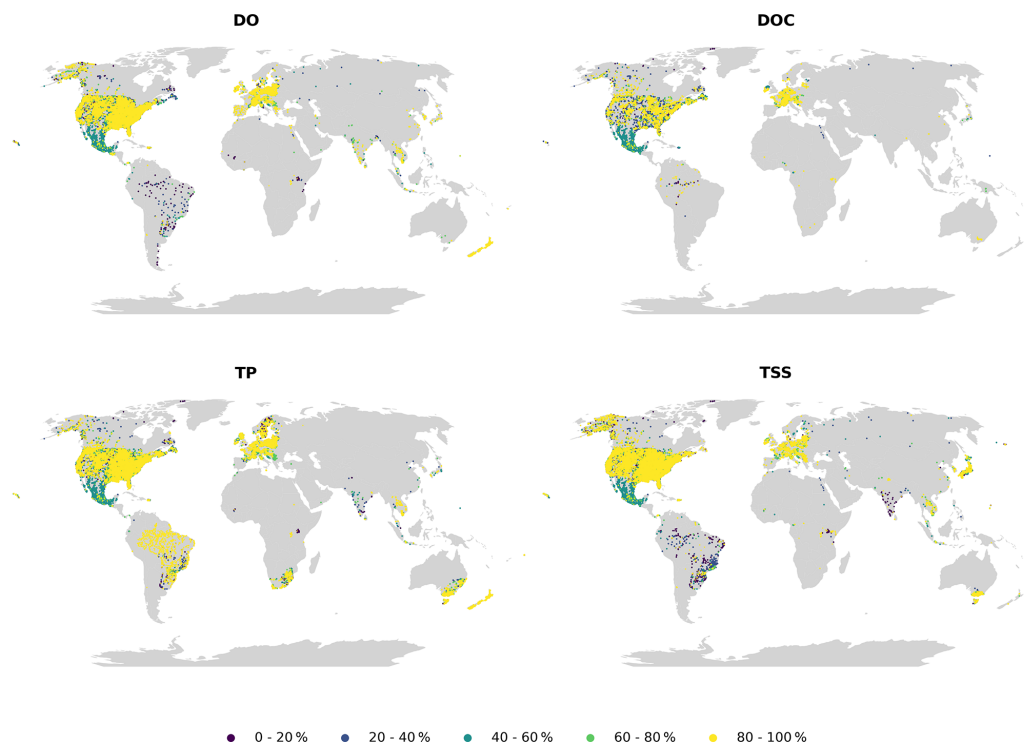


Figure 5. Monthly availability for dissolved oxygen (DO), dissolved organic carbon (DOC), total phosphorus (TP) and total suspended solids (TSSs).

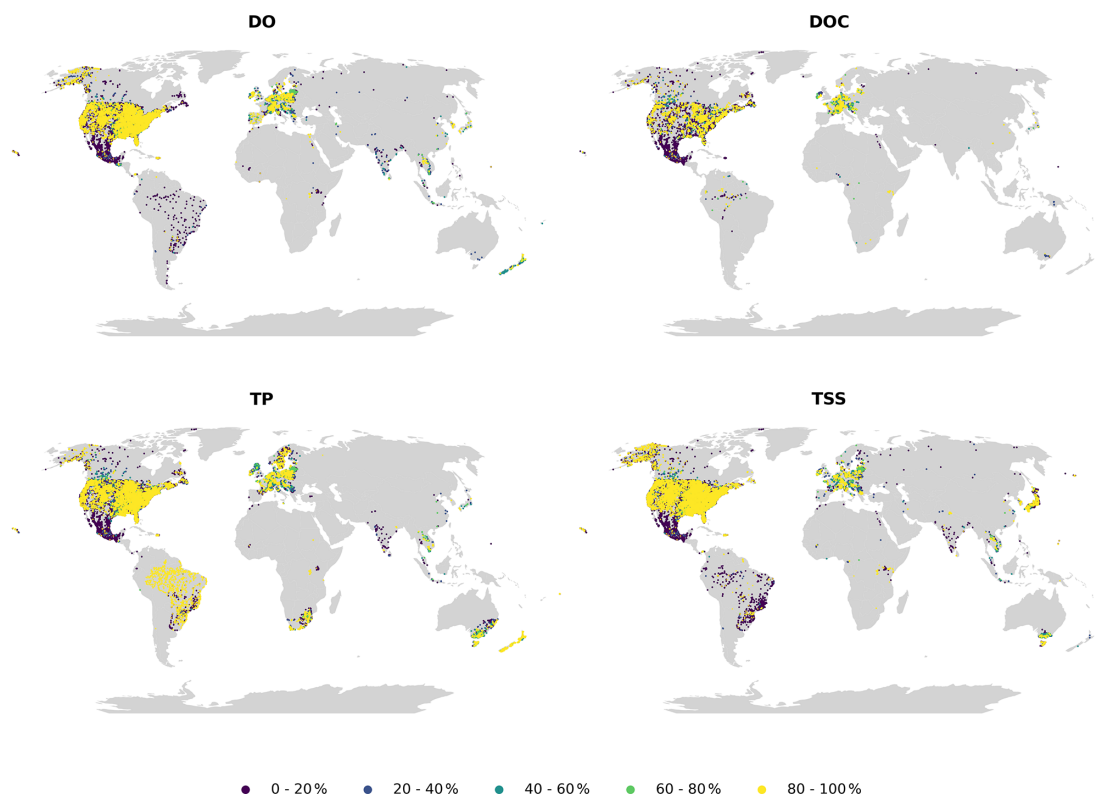


Figure 6. Monthly continuity for dissolved oxygen (DO), dissolved organic carbon (DOC), total phosphorus (TP) and total suspended solids (TSSs).

source data compilation that can be customized for other data sources as well. The current version of GRQA is mainly focused on different forms of the main nutrients (N and P) and carbon compounds, although GEMStat, Waterbase and WQP also have many other types of parameters that are used as water quality indicators (heavy metals, pesticides, etc.). Other researchers are able to make additions and customize the dataset to their needs for parameter-specific studies using the scripts published with GRQA.

Updates and additions by the hydrological community are encouraged to further develop GRQA. As it stands, GRQA is a set of well-structured CSV files rather than a database able to be queried. We intend to add a Jupyter Notebook example of loading and processing the CSV files to the GRQA GitHub repository. We include an extensive data catalog with graphs and maps for temporal and spatial coverage of every variable in the Supplement. This should help potential users to get a better overview of the data before downloading it. Converting the files into a database would also greatly improve data management and make extending GRQA easier in the future. In the case of a relational database, the schema recommended by Plana et al. (2019) could be followed. We are also considering the addition of an online dashboard for data visualization and download, similar to that of GEMStat or WQP. A versioning system along with a metadata validation strategy similar to Welty et al. (2020) could be implemented to ensure metadata quality.

Future work could also include the development of a dataset for catchment characteristics in order to better study how water quality in rivers and streams is affected by land use changes in their catchments. The CAMELS dataset (Addor et al., 2017) and its regional implementations (Chagas et al., 2020; Coxon et al., 2020) can be used as an example. In addition, interactions between water quality and streamflow can be further studied by linking water quality observations to streamflow data from the Global Streamflow Indices and Metadata Archive (GSIM) (Do et al., 2018).

Appendix A: Figures and tables in appendices

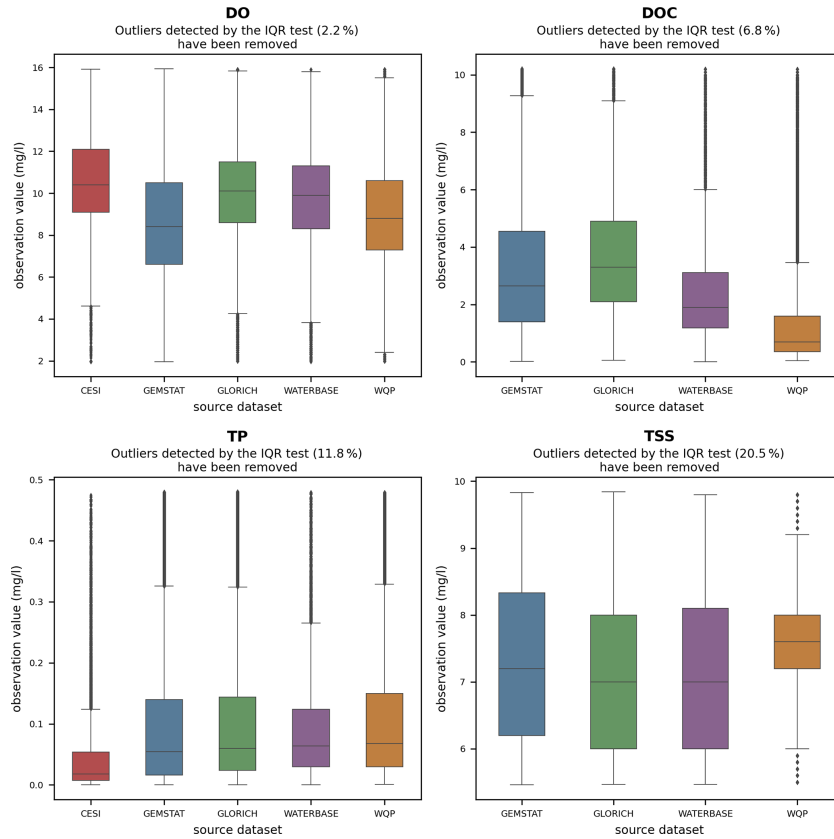


Figure A1. Box plot of observation values for dissolved oxygen (DO), dissolved organic carbon (DOC), total phosphorus (TP) and total suspended solids (TSSs). Outliers determined by the IQR test (Table 7) are not shown on the plot.

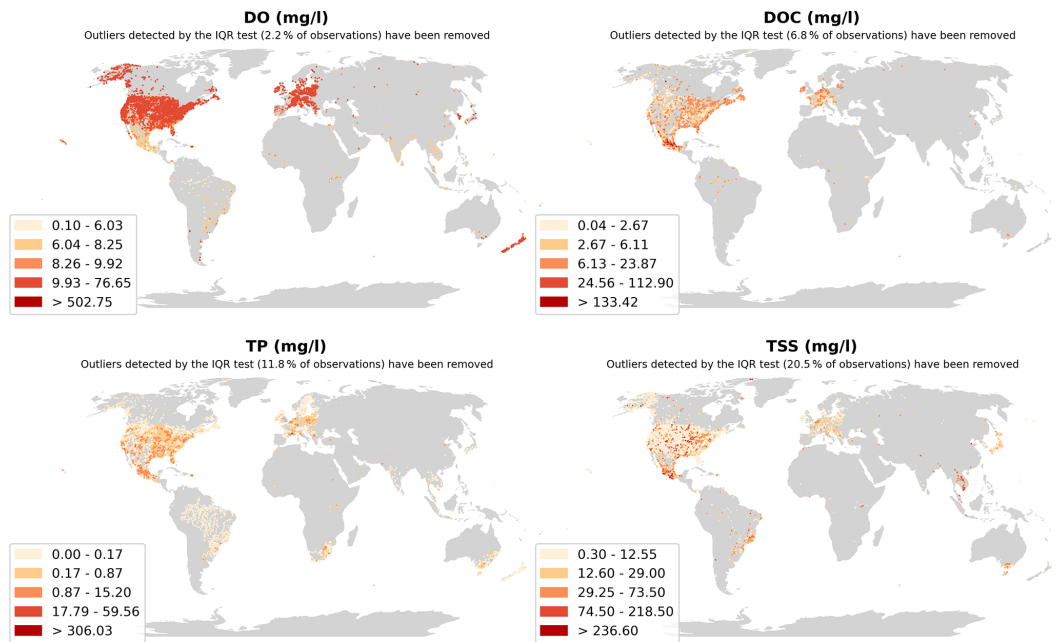


Figure A2. Spatial distribution of yearly median observation values for dissolved oxygen (DO), dissolved organic carbon (DOC), total phosphorus (TP) and total suspended solids (TSSs). Outliers determined by the IQR test are not shown on the plot.

Table A1. Conversion procedures of source data units and chemical forms into their corresponding GRQA versions for all parameters.

Parameter code	Source	Form	Source form	Unit	Source unit	Divisor	Multiplier	Conversion constant
TAN	CESI	N	NH ₃	mg L ⁻¹	mg L ⁻¹	17.031	14.007	0.822441
NO3N	CESI	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
NO2N	CESI	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
TN	CESI	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
TDN	CESI	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
DO	CESI	O ₂	O ₂	mg L ⁻¹	mg L ⁻¹	1	1	1
pH	CESI			pH	pH units	1	1	1
TP	CESI	P	P	mg L ⁻¹	mg L ⁻¹	1	1	1
TDP	CESI	P	P	mg L ⁻¹	mg L ⁻¹	1	1	1
TEMP	CESI			°C	°C	1	1	1
DC	GEMStat	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
DIC	GEMStat	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
DOC	GEMStat	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
POC	GEMStat	C	C	mg L ⁻¹	µg g ⁻¹	1	1	1
POC	GEMStat	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
TC	GEMStat	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
TIC	GEMStat	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
TOC	GEMStat	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
DKN	GEMStat	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
DON	GEMStat	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
NH4N	GEMStat	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
NH4N	GEMStat	N	NH ₄	mg L ⁻¹	mg L ⁻¹ NH ₄	18.039	14.007	0.776484
NH4N	GEMStat	N	N	mg L ⁻¹	µg L ⁻¹	1000	1	0.001
NO2N	GEMStat	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
NO2N	GEMStat	N	NO ₂	mg L ⁻¹	mg L ⁻¹ NO ₂	46.005	14.007	0.304467
NO2N	GEMStat	N	N	mg L ⁻¹	µg L ⁻¹	1000	1	0.001
NO3N	GEMStat	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
NO3N	GEMStat	N	NO ₃	mg L ⁻¹	mg L ⁻¹ NO ₃	62.004	14.007	0.225905
NO3N	GEMStat	N	N	mg L ⁻¹	µg L ⁻¹	1000	1	0.001
PN	GEMStat	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
PON	GEMStat	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
PON	GEMStat	N	N	mg L ⁻¹	µg g ⁻¹	1	1	1
TDN	GEMStat	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
TKN	GEMStat	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
TN	GEMStat	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
TON	GEMStat	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
DO	GEMStat	O ₂	O ₂	mg L ⁻¹	mg L ⁻¹	1	1	1
DOSAT	GEMStat			%	%	1	1	1
BOD	GEMStat	O ₂	O ₂	mg L ⁻¹	mg L ⁻¹	1	1	1
COD	GEMStat	O ₂	O ₂	mg L ⁻¹	mg L ⁻¹	1	1	1
pH	GEMStat			pH	–	1	1	1
DIP	GEMStat	P	P	mg L ⁻¹	mg L ⁻¹	1	1	1
TDP	GEMStat	P	P	mg L ⁻¹	mg L ⁻¹	1	1	1
TDP	GEMStat	P	P	mg L ⁻¹	µg L ⁻¹	1000	1	0.001
TIP	GEMStat	P	P	mg L ⁻¹	mg L ⁻¹	1	1	1
TP	GEMStat	P	P	mg L ⁻¹	mg L ⁻¹	1	1	1
TP	GEMStat	P	P	mg L ⁻¹	µg L ⁻¹	1000	1	0.001
TPP	GEMStat	P	P	mg L ⁻¹	µg g ⁻¹	1	1	1
TPP	GEMStat	P	P	mg L ⁻¹	mg L ⁻¹	1	1	1
TSS	GEMStat			mg L ⁻¹	mg L ⁻¹	1	1	1
TEMP	GEMStat			°C	°C	1	1	1
TEMP	GLORICH			°C	°C	1	1	1

Table A1. Continued.

Parameter code	Source	Form	Source form	Unit	Source unit	Divisor	Multiplier	Conversion constant
pH	GLORICH			pH		1	1	1
DO	GLORICH	O ₂	O ₂	mg L ⁻¹	mg O ₂ L ⁻¹	1	1	1
DOSAT	GLORICH			%	%	1	1	1
TSS	GLORICH			mg L ⁻¹	mg L ⁻¹	1	1	1
TC	GLORICH	C	C	mg L ⁻¹	μmol L ⁻¹	1000	12.011	0.012011
TIC	GLORICH	C	C	mg L ⁻¹	μmol L ⁻¹	1000	12.011	0.012011
DIC	GLORICH	C	C	mg L ⁻¹	μmol L ⁻¹	1000	12.011	0.012011
PIC	GLORICH	C	C	mg L ⁻¹	μmol L ⁻¹	1000	12.011	0.012011
TOC	GLORICH	C	C	mg L ⁻¹	μmol L ⁻¹	1000	12.011	0.012011
DOC	GLORICH	C	C	mg L ⁻¹	μmol L ⁻¹	1000	12.011	0.012011
POC	GLORICH	C	C	mg L ⁻¹	μmol L ⁻¹	1000	12.011	0.012011
TN	GLORICH	N	N	mg L ⁻¹	μmol L ⁻¹	1000	14.007	0.014007
TDN	GLORICH	N	N	mg L ⁻¹	μmol L ⁻¹	1000	14.007	0.014007
PN	GLORICH	N	N	mg L ⁻¹	μmol L ⁻¹	1000	14.007	0.014007
TIN	GLORICH	N	N	mg L ⁻¹	μmol L ⁻¹	1000	14.007	0.014007
DIN	GLORICH	N	N	mg L ⁻¹	μmol L ⁻¹	1000	14.007	0.014007
TON	GLORICH	N	N	mg L ⁻¹	μmol L ⁻¹	1000	14.007	0.014007
DON	GLORICH	N	N	mg L ⁻¹	μmol L ⁻¹	1000	14.007	0.014007
PON	GLORICH	N	N	mg L ⁻¹	μmol L ⁻¹	1000	14.007	0.014007
TKN	GLORICH	N	N	mg L ⁻¹	μmol L ⁻¹	1000	14.007	0.014007
DKN	GLORICH	N	N	mg L ⁻¹	μmol L ⁻¹	1000	14.007	0.014007
NO3N	GLORICH	N	NO ₃	mg L ⁻¹	μmol L ⁻¹	1000	0.225905	0.000226
NO2N	GLORICH	N	NO ₂	mg L ⁻¹	μmol L ⁻¹	1000	0.304467	0.000304
NH4N	GLORICH	N	NH ₄	mg L ⁻¹	μmol L ⁻¹	1000	0.776484	0.000776
TP	GLORICH	P	P	mg L ⁻¹	μmol L ⁻¹	1000	30.973	0.030973
TDP	GLORICH	P	P	mg L ⁻¹	μmol L ⁻¹	1000	30.973	0.030973
TPP	GLORICH	P	P	mg L ⁻¹	μmol L ⁻¹	1000	30.973	0.030973
TIP	GLORICH	P	P	mg L ⁻¹	μmol L ⁻¹	1000	30.973	0.030973
DIP	GLORICH	P	P	mg L ⁻¹	μmol L ⁻¹	1000	30.973	0.030973
NO3N	Waterbase	N	NO ₃	mg L ⁻¹	mg NO ₃ L ⁻¹	62.004	14.007	0.225905
NO2N	Waterbase	N	NO ₂	mg L ⁻¹	mg NO ₂ L ⁻¹	46.005	14.007	0.304467
NH4N	Waterbase	N	NH ₄	mg L ⁻¹	mg NH ₄ L ⁻¹	18.039	14.007	0.776484
NH4N	Waterbase	N	NH ₃	mg L ⁻¹	mg NH ₃ L ⁻¹	17.031	14.007	0.822441
NH3N	Waterbase	N	NH ₃	mg L ⁻¹	mg NH ₃ L ⁻¹	17.031	14.007	0.822441
NH3N	Waterbase	N	N	mg L ⁻¹	μg L ⁻¹	1000	1	0.001
TP	Waterbase	P	P	mg L ⁻¹	mg P L ⁻¹	1	1	1
TSS	Waterbase			mg L ⁻¹	mg L ⁻¹	1	1	1
TEMP	Waterbase			°C	Celsius	1	1	1
DOSAT	Waterbase			%	%	1	1	1
DO	Waterbase	O ₂	O ₂	mg L ⁻¹	mg L ⁻¹	1	1	1
DO	Waterbase	O ₂	O ₂	mg L ⁻¹	mg O ₂ L ⁻¹	1	1	1
BOD5	Waterbase	O ₂	O ₂	mg L ⁻¹	mg O ₂ L ⁻¹	1	1	1
BOD7	Waterbase	O ₂	O ₂	mg L ⁻¹	mg O ₂ L ⁻¹	1	1	1
CODCr	Waterbase	O ₂	O ₂	mg L ⁻¹	mg O ₂ L ⁻¹	1	1	1
CODMn	Waterbase	O ₂	O ₂	mg L ⁻¹	mg O ₂ L ⁻¹	1	1	1
DOC	Waterbase	C	C	mg L ⁻¹	mg C L ⁻¹	1	1	1
DOC	Waterbase	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
TOC	Waterbase	C	C	mg L ⁻¹	mg C L ⁻¹	1	1	1
TOC	Waterbase	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
pH	Waterbase			pH		1	1	1
TKN	Waterbase	N	N	mg L ⁻¹	mg N L ⁻¹	1	1	1
TKN	Waterbase	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1

Table A1. Continued.

Parameter code	Source	Form	Source form	Unit	Source unit	Divisor	Multiplier	Conversion constant
TON	Waterbase	N	N	mg L ⁻¹	mg N L ⁻¹	1	1	1
PON	Waterbase	N	N	mg L ⁻¹	mg N L ⁻¹	1	1	1
TIN	Waterbase	N	N	mg L ⁻¹	mg N L ⁻¹	1	1	1
TN	Waterbase	N	N	mg L ⁻¹	mg N L ⁻¹	1	1	1
PC	WQP	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
DC	WQP	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
TC	WQP	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
DO	WQP	O ₂	O ₂	mg L ⁻¹	mg L ⁻¹	1	1	1
DOSAT	WQP			%	% saturation	1	1	1
PIC	WQP	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
DIC	WQP	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
TIC	WQP	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
TAN	WQP	N	N	mg L ⁻¹	mg L ⁻¹ as N	1	1	1
TAN	WQP	N	N	mg L ⁻¹	mg L ⁻¹ as N	1	1	1
DIN	WQP	N	N	mg L ⁻¹	mg L ⁻¹ as N	1	1	1
TIN	WQP	N	N	mg L ⁻¹	mg L ⁻¹ as N	1	1	1
NO3N	WQP	N	N	mg L ⁻¹	mg L ⁻¹ as N	1	1	1
NO3N	WQP	N	N	mg L ⁻¹	mg L ⁻¹ as N	1	1	1
NO2N	WQP	N	N	mg L ⁻¹	mg L ⁻¹ as N	1	1	1
NO2N	WQP	N	N	mg L ⁻¹	mg L ⁻¹ as N	1	1	1
PON	WQP	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
DON	WQP	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
TON	WQP	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
POP	WQP	P	P	mg L ⁻¹	mg L ⁻¹ as P	1	1	1
DOP	WQP	P	P	mg L ⁻¹	mg L ⁻¹ as P	1	1	1
TOP	WQP	P	P	mg L ⁻¹	mg L ⁻¹ as P	1	1	1
PN	WQP	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
TPP	WQP	P	P	mg L ⁻¹	mg L ⁻¹ as P	1	1	1
TDP	WQP	P	P	mg L ⁻¹	mg L ⁻¹ as P	1	1	1
TP	WQP	P	P	mg L ⁻¹	mg L ⁻¹ as P	1	1	1
TP	WQP	P	P	mg L ⁻¹	mg L ⁻¹ as P	1	1	1
TN	WQP	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
TDN	WQP	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
TN	WQP	N	N	mg L ⁻¹	mg L ⁻¹	1	1	1
POC	WQP	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
DOC	WQP	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
TOC	WQP	C	C	mg L ⁻¹	mg L ⁻¹	1	1	1
BOD5	WQP	O ₂	O ₂	mg L ⁻¹	mg L ⁻¹	1	1	1
BOD5	WQP	O ₂	O ₂	mg L ⁻¹	mg L ⁻¹	1	1	1
pH	WQP			pH	Standard units	1	1	1
TSS	WQP			mg L ⁻¹	mg L ⁻¹	1	1	1
TEMP	WQP			°C	°C	1	1	1

Supplement. The supplement related to this article is available online at: <https://doi.org/10.5194/essd-13-5483-2021-supplement>.

Author contributions. HV, AK and EU conceived the idea of the study. GA and EU aided with establishing a theoretical background for the study. HV designed the code and carried out the data processing with contributions from AK. HV, GA and LS developed the observation deduplication and outlier detection strategy. HV prepared the manuscript with contributions from all co-authors.

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

Disclaimer. Publisher's note: Copernicus Publications remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Acknowledgements. This research was funded by Mobilitas+ program grant no. MOBERC34, Marie Skłodowska-Curie Actions individual fellowship under the Horizon 2020 Programme grant agreement number 795625, grant PRG352 from the Estonian Research Council, NUTIKAS program, and the Dora Plus PhD student mobility scholarship number 36.9-6.1/1124 of the Archimedes foundation and European Regional Development Fund (EcolChange Centre of Excellence). Holger Virro is also thankful for technical support from the Yale Center for Research Computing support and the High Performance Computing Center of the University of Tartu.

Financial support. This research has been supported by the Sihtasutus Archimedes (grant no. 36.9-6.1/1124), the Eesti Teadusagentuur (grant no. MOBERC34), the H2020 Marie Skłodowska-Curie Actions (grant no. 795625) and the Eesti Teadusagentuur (grant no. PRG352).

Review statement. This paper was edited by Birgit Heim and reviewed by two anonymous referees.

References

- Abbaspour, K. C., Rouholahnejad, E., Vaghefi, S., Srinivasan, R., Yang, H., and Kløve, B.: A continental-scale hydrology and water quality model for Europe: Calibration and uncertainty of a high-resolution large-scale SWAT model, *J. Hydrol.*, 524, 733–752, <https://doi.org/10.1016/j.jhydrol.2015.03.027>, 2015.
- Addor, N., Newman, A. J., Mizukami, N., and Clark, M. P.: The CAMELS data set: catchment attributes and meteorology for large-sample studies, *Hydrol. Earth Syst. Sci.*, 21, 5293–5313, <https://doi.org/10.5194/hess-21-5293-2017>, 2017.
- Archfield, S. A., Clark, M., Arheimer, B., Hay, L. E., McMillan, H., Kiang, J. E., Seibert, J., Hakala, K., Bock, A., Wagener, T., Farmer, W. H., Andréassian, V., Attinger, S., Viglione, A., Knight, R., Markstrom, S., and Over, T.: Accelerating advances in continental domain hydrologic modeling, *Water Resour. Res.*, 51, 10078–10091, 2015.
- Beck, H. E., De Roo, A., and van Dijk, A. I.: Global maps of stream-flow characteristics based on observations from several thousand catchments, *J. Hydrometeorol.*, 16, 1478–1501, 2015.
- Berndt, D. J. and Clifford, J.: Using dynamic time warping to find patterns in time series, in: KDD workshop, Seattle, WA, USA, 26 April 1994, 10, 359–370, available at: <https://www.aaai.org/Papers/Workshops/1994/WS-94-03/WS94-03-031.pdf> (last access: 27 January 2021), 1994.
- Bierkens, M. F.: Global hydrology 2015: State, trends, and directions, *Water Resour. Res.*, 51, 4923–4947, 2015.
- Birant, D. and Kut, A.: ST-DBSCAN: An algorithm for clustering spatial-temporal data, *Data Knowl. Eng.*, 60, 208–221, 2007.
- Blöschl, G., Bierkens, M. F., Chambel, A., et al.: Twenty-three unsolved problems in hydrology (UPH) – a community perspective, *Hydrolog. Sci. J.*, 64, 1141–1158, 2019.
- Börker, J., Hartmann, J., Amann, T., Romero-Mujalli, G., Moosdorf, N., and Jenkins, C.: Chemical river data from drained loess areas, PANGAEA [data set], <https://doi.org/10.1594/PANGAEA.915784>, 2020.
- Box, G. E. and Cox, D. R.: An analysis of transformations, *J. Roy. Stat. Soc. B-Met.*, 26, 211–243, 1964.
- Caraco, N. F. and Cole, J. J.: Human impact on nitrate export: an analysis using major world rivers, *Ambio*, 28, 167–170, 1999.
- Castrillo, M. and García, Á. L.: Estimation of high frequency nutrient concentrations from water quality surrogates using machine learning methods, *Water Res.*, 172, 115490, <https://doi.org/10.1016/j.watres.2020.115490>, 2020.
- Chagas, V. B. P., Chaffe, P. L. B., Addor, N., Fan, F. M., Fleischmann, A. S., Paiva, R. C. D., and Siqueira, V. A.: CAMELS-BR: hydrometeorological time series and landscape attributes for 897 catchments in Brazil, *Earth Syst. Sci. Data*, 12, 2075–2096, <https://doi.org/10.5194/essd-12-2075-2020>, 2020.
- Chau, K.-w.: A review on integration of artificial intelligence into water quality modelling, *Mar. Pollut. Bull.*, 52, 726–733, 2006.
- Chen, J. and Quan, W.: Using Landsat/TM imagery to estimate nitrogen and phosphorus concentration in Taihu Lake, China, *IEEE J. Sel. Top. Appl.*, 5, 273–280, 2011.
- Chen, K., Chen, H., Zhou, C., Huang, Y., Qi, X., Shen, R., Liu, F., Zuo, M., Zou, X., Wang, J., et al.: Comparative analysis of surface water quality prediction performance and identification of key water parameters using different machine learning models based on big data, *Water Res.*, 171, 115454, <https://doi.org/10.1016/j.watres.2019.115454>, 2020.
- Choubin, B., Darabi, H., Rahmati, O., Sajedi-Hosseini, F., and Kløve, B.: River suspended sediment modelling using the CART model: a comparative study of machine learning techniques, *Sci. Total Environ.*, 615, 272–281, 2018.
- Coxon, G., Addor, N., Bloomfield, J. P., Freer, J., Fry, M., Hannaford, J., Howden, N. J. K., Lane, R., Lewis, M., Robinson, E. L., Wagener, T., and Woods, R.: CAMELS-GB: hydrometeorological time series and landscape attributes for 671 catchments in Great Britain, *Earth Syst. Sci. Data*, 12, 2459–2483, <https://doi.org/10.5194/essd-12-2459-2020>, 2020.
- Crochemore, L., Isberg, K., Pimentel, R., Pineda, L., Hasan, A., and Arheimer, B.: Lessons learnt from checking the quality of openly

- accessible river flow data worldwide, *Hydrolog. Sci. J.*, 65, 699–711, <https://doi.org/10.1080/02626667.2019.1659509>, 2019.
- Desmit, X., Thieu, V., Billen, G., Campuzano, F., Dulière, V., Garnier, J., Lassaletta, L., Ménesguen, A., Neves, R., Pinto, L., Silvestre M., Sobrinho, J. L., and Lacroix, G.: Reducing marine eutrophication may require a paradigmatic change, *Sci. Total Environ.*, 635, 1444–1466, 2018.
- Do, H. X., Gudmundsson, L., Leonard, M., and 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, <https://doi.org/10.5194/essd-10-765-2018>, 2018.
- Enderlein, R., Enderlein, R., and Williams, W.: Chapter 2: Water Quality Requirements, in: *Water Pollution Control – A Guide to the Use of Water, Quality Management Principles*, edited by: Helmer, R. and Hespanhol, I., WHO/UNEP, 1996.
- Environment and Climate Change Canada: Water quality in Canadian rivers, available at: <https://open.canada.ca/data/en/dataset/55cc50dc-feb3-46d1-b40f-09254f3c00c5>, last access: 16 November 2020.
- European Environment Agency: Waterbase – Water Quality ICM, available at: <https://www.eea.europa.eu/data-and-maps/data/waterbase-water-quality-icm>, last access: 16 November 2020.
- Evans, C., Monteith, D., and Cooper, D.: Long-term increases in surface water dissolved organic carbon: observations, possible causes and environmental impacts, *Environ. Pollut.*, 137, 55–71, 2005.
- Foley, J. A., Ramankutty, N., Brauman, K. A., Cassidy, E. S., Gerber, J. S., Johnston, M., Mueller, N. D., O’Connell, C., Ray, D. K., West, P. C., Balzer, C., Bennett, E. M., Carpenter, S. R., Hill, J., Monfreda, C., Polasky, S., Rockström, J., Sheehan, J., Siebert, S., Tilman, D., and Zaks, D. P. M.: Solutions for a cultivated planet, *Nature*, 478, 337–342, 2011.
- Gao, Y., Merz, C., Lischeid, G., and Schneider, M.: A review on missing hydrological data processing, *Environ. Earth Sci.*, 77, 47, <https://doi.org/10.1007/s12665-018-7228-6>, 2018.
- Gudivada, V., Apon, A., and Ding, J.: Data quality considerations for big data and machine learning: Going beyond data cleaning and transformations, *International Journal on Advances in Software*, 10, 1–20, 2017.
- Gudmundsson, L. and Seneviratne, S. I.: Towards observation-based gridded runoff estimates for Europe, *Hydrol. Earth Syst. Sci.*, 19, 2859–2879, <https://doi.org/10.5194/hess-19-2859-2015>, 2015.
- Harrigan, S., Zsoter, E., Alfieri, L., Prudhomme, C., Salamon, P., Wetterhall, F., Barnard, C., Cloke, H., and Pappenberger, F.: GloFAS-ERA5 operational global river discharge reanalysis 1979–present, *Earth Syst. Sci. Data*, 12, 2043–2060, <https://doi.org/10.5194/essd-12-2043-2020>, 2020.
- Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M., Haldane, A., del Río, J. F., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K., Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C., Oliphant, and Travis, E.: Array programming with NumPy, *Nature*, 585, 357–362, 2020.
- Hartmann, J., Lauerwald, R., and Moosdorf, N.: A Brief Overview of the GLObal River Chemistry Database, *GLORICH, Proced. Earth Plan. Sc.*, 10, 23–27, <https://doi.org/10.1016/J.PROEPS.2014.08.005>, 2014.
- Hartmann, J., Lauerwald, R., and Moosdorf, N.: GLORICH-Global river chemistry database, PANGAEA [data set], <https://doi.org/10.1594/PANGAEA.902360>, 2019.
- He, B., Kanae, S., Oki, T., Hirabayashi, Y., Yamashiki, Y., and Takara, K.: Assessment of global nitrogen pollution in rivers using an integrated biogeochemical modeling framework, *Water Res.*, 45, 2573–2586, 2011.
- Helsel, D. R.: Advantages of nonparametric procedures for analysis of water quality data, *Hydrolog. Sci. J.*, 32, 179–190, 1987.
- Hirsch, R. M., Slack, J. R., and Smith, R. A.: Techniques of trend analysis for monthly water quality data, *Water Resour. Res.*, 18, 107–121, 1982.
- Hope, D., Billett, M., and Cresser, M.: A review of the export of carbon in river water: fluxes and processes, *Environ. Pollut.*, 84, 301–324, 1994.
- Hrachowitz, M., Savenije, H., Blöschl, G., McDonnell, J., Sivalalan, M., Pomeroy, J., Arheimer, B., Blume, T., Clark, M., Ehret, U., Fenicia, F., Freer, J. E., Gelfan, A., Gupta, H. V., Hughes, D. A., Hut, R. W., Montanari, A., Pande, S., Tetzlaff, D., Troch, P. A., Uhlenbrook, S., Wagener, T., Winsemius, H. C., Woods, R. A., Zehe, E., and Cudennec, C.: A decade of Predictions in Ungauged Basins (PUB) a review, *Hydrolog. Sci. J.*, 58, 1198–1255, 2013.
- Hughes, A. O., Tanner, C. C., McKergow, L. A., and Sukias, J. P.: Unrestricted dairy cattle grazing of a pastoral headwater wetland and its effect on water quality, *Agr. Water Manage.*, 165, 72–81, 2016.
- Hutton, C., Wagener, T., Freer, J., Han, D., Duffy, C., and Arheimer, B.: Most computational hydrology is not reproducible, so is it really science?, *Water Resour. Res.*, 52, 7548–7555, 2016.
- Jones, A. S., Stevens, D. K., Horsburgh, J. S., and Mesner, N. O.: Surrogate Measures for Providing High Frequency Estimates of Total Suspended Solids and Total Phosphorus Concentrations, *J. Am. Water Resour. As.*, 47, 239–253, 2011.
- Jordahl, K., den Bossche, J. V., Fleischmann, M., Wasserman, J., McBride, J., Gerard, J., Tratner, J., Perry, M., Badaracco, A. G., Farmer, C., Hjelle, G. A., Snow, A. D., Cochran, M., Gillies, S., Culbertson, L., Bartos, M., Eubank, N., maxalbert, Bilogur, A., Rey, S., Ren, C., Arribas-Bel, D., Wasser, L., Wolf, L. J., Journois, M., Wilson, J., Greenhall, A., Holdgraf, C., Filipe, and Leblanc, F.: *geopandas/geopandas: v0.8.1*, Zenodo [code], <https://doi.org/10.5281/zenodo.3946761>, 2020.
- Khan, K., Rehman, S. U., Aziz, K., Fong, S., and Sarasvady, S.: DB-SCAN: Past, present and future, *The Fifth International Conference on the Applications of Digital Information and Web Technologies (ICADIWT 2014)*, 232–238, 2014.
- Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A. K., Hochreiter, S., and Nearing, G. S.: Toward improved predictions in ungauged basins: Exploiting the power of machine learning, *Water Resour. Res.*, 55, 11344–11354, 2019.
- Krysanova, V., Müller-Wohlfeil, D.-I., and Becker, A.: Development and test of a spatially distributed hydrological/water quality model for mesoscale watersheds, *Ecol. Model.*, 106, 261–289, 1998.
- Leon, L., Soulis, E., Kouwen, N., and Farquhar, G.: Nonpoint source pollution: a distributed water quality modeling approach, *Water Res.*, 35, 997–1007, 2001.

- Marzadri, A., Amatulli, G., Tonina, D., Bellin, A., Shen, L. Q., Allen, G. H., and Raymond, P. A.: Global riverine nitrous oxide emissions: the role of small streams and large rivers, *Sci. Total Environ.*, 776, 145148, <https://doi.org/10.1016/j.scitotenv.2021.145148>, 2021.
- McKinney, W.: Data structures for statistical computing in python, in: Proceedings of the 9th Python in Science Conference, Austin, TX, 28 June and 3 July 2010, 445, 51–56, <https://doi.org/10.25080/Majora-92bf1922-012>, 2010.
- McMillan, H., Krueger, T., and Freer, J.: Benchmarking observational uncertainties for hydrology: rainfall, river discharge and water quality, *Hydrol. Process.*, 26, 4078–4111, 2012.
- Meals, D. W., Dressing, S. A., and Davenport, T. E.: Lag time in water quality response to best management practices: A review, *J. Environ. Qual.*, 39, 85–96, 2010.
- Mount, N. J., Maier, H. R., Toth, E., Elshorbagy, A., Solomatine, D., Chang, F.-J., and Abrahart, R.: Data-driven modelling approaches for socio-hydrology: opportunities and challenges within the Panta Rhei Science Plan, *Hydrolog. Sci. J.*, 61, 1192–1208, 2016.
- Mueller, N. D., Gerber, J. S., Johnston, M., Ray, D. K., Ramankutty, N., and Foley, J. A.: Closing yield gaps through nutrient and water management, *Nature*, 490, 254–257, 2012.
- Neukermans, G., Ruddick, K., Loisel, H., and Roose, P.: Optimization and quality control of suspended particulate matter concentration measurement using turbidity measurements, *Limnol. Oceanogr.-Meth.*, 10, 1011–1023, <https://doi.org/10.4319/lom.2012.10.1011>, 2012.
- Olmanson, L. G., Brezonik, P. L., and Bauer, M. E.: Airborne hyperspectral remote sensing to assess spatial distribution of water quality characteristics in large rivers: The Mississippi River and its tributaries in Minnesota, *Remote Sens. Environ.*, 130, 254–265, 2013.
- Ouali, D., Chebana, F., and Ouarda, T. B.: Fully nonlinear statistical and machine-learning approaches for hydrological frequency estimation at ungauged sites, *J. Adv. Model. Earth Sy.*, 9, 1292–1306, 2017.
- Ouyang, W., Yang, W., Tysklind, M., Xu, Y., Lin, C., Gao, X., and Hao, Z.: Using river sediments to analyze the driving force difference for non-point source pollution dynamics between two scales of watersheds, *Water Res.*, 139, 311–320, 2018.
- Papacharalampous, G., Tyralis, H., and Koutsoyiannis, D.: Comparison of stochastic and machine learning methods for multi-step ahead forecasting of hydrological processes, *Stoch. Env. Res. Risk A.*, 33, 481–514, 2019.
- Parimala, M., Lopez, D., and Senthilkumar, N.: A survey on density based clustering algorithms for mining large spatial databases, *International Journal of Advanced Science and Technology*, 31, 59–66, 2011.
- Parmar, K. S. and Bhardwaj, R.: Water quality management using statistical analysis and time-series prediction model, *Applied Water Science*, 4, 425–434, 2014.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al.: Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, É.: Scikit-learn: Machine learning in Python, *J. Mach. Learn. Res.*, 12, 2825–2830, 2011.
- Pellerin, B. A., Stauffer, B. A., Young, D. A., Sullivan, D. J., Bricker, S. B., Walbridge, M. R., Clyde Jr., G. A., and Shaw, D. M.: Emerging tools for continuous nutrient monitoring networks: Sensors advancing science and water resources protection, *J. Am. Water Resour. As.*, 52, 993–1008, 2016.
- Plana, Q., Alferes, J., Fuks, K., Kraft, T., Maruéjols, T., Torfs, E., and Vanrolleghem, P. A.: Towards a water quality database for raw and validated data with emphasis on structured metadata, *Water Qual. Res. J.*, 54, 1–9, 2019.
- Radwan, M., Willems, P., El-Sadek, A., and Berlamont, J.: Modelling of dissolved oxygen and biochemical oxygen demand in river water using a detailed and a simplified model, *International Journal of River Basin Management*, 1, 97–103, 2003.
- Read, E. K., Carr, L., De Cicco, L., Dugan, H. A., Hanson, P. C., Hart, J. A., Kreft, J., Read, J. S., and Winslow, L. A.: Water quality data for national-scale aquatic research: The Water Quality Portal, *Water Resour. Res.*, 53, 1735–1745, 2017.
- Ren, H., Cromwell, E., Kravitz, B., and Chen, X.: Using Deep Learning to Fill Spatio-Temporal Data Gaps in Hydrological Monitoring Networks, *Hydrol. Earth Syst. Sci. Discuss.* [preprint], <https://doi.org/10.5194/hess-2019-196>, in review, 2019.
- Shen, C., Laloy, E., Elshorbagy, A., Albert, A., Bales, J., Chang, F.-J., Ganguly, S., Hsu, K.-L., Kifer, D., Fang, Z., Fang, K., Li, D., Li, X., and Tsai, W.-P.: HESS Opinions: Incubating deep-learning-powered hydrologic science advances as a community, *Hydrol. Earth Syst. Sci.*, 22, 5639–5656, <https://doi.org/10.5194/hess-22-5639-2018>, 2018.
- Shen, L. Q., Amatulli, G., Sethi, T., Raymond, P., and Domisch, S.: Estimating nitrogen and phosphorus concentrations in streams and rivers, within a machine learning framework, *Scientific Data*, 7, 161, <https://doi.org/10.1038/s41597-020-0478-7>, 2020.
- Singh, K. P., Basant, A., Malik, A., and Jain, G.: Artificial neural network modeling of the river water quality – a case study, *Ecol. Model.*, 220, 888–895, 2009.
- Sinha, E., Michalak, A., Calvin, K. V., and Lawrence, P. J.: Societal decisions about climate mitigation will have dramatic impacts on eutrophication in the 21 st century, *Nat. Commun.*, 10, 939, <https://doi.org/10.1038/s41467-019-08884-w>, 2019.
- Snow, A. D., Whitaker, J., Cochran, M., den Bossche, J. V., Mayo, C., de Kloe, J., Karney, C., Ouzounoudis, G., Dearing, J., Lostis, G., Heitor, Filipe, May, R., Itkin, M., Couwenberg, B., Berardinelli, G., Badger, T. G., Eubank, N., Dunphy, M., Brett, M., Raspaud, M., da Costa, M. A., Evers, K., Ranalli, J., de Maeyer, J., Popov, E., Gohlke, C., Willoughby, C., Barker, C., and Wiedemann, B. M.: pyproj4/pyproj: 2.6.1 Release, Zenodo [code], <https://doi.org/10.5281/zenodo.3783866>, 2020.
- Sprague, L. A., Oelsner, G. P., and Argue, D. M.: Challenges with secondary use of multi-source water-quality data in the United States, *Water Res.*, 110, 252–261, 2017.
- Stagge, J. H., Rosenberg, D. E., Abdallah, A. M., Akbar, H., Attallah, N. A., and James, R.: Assessing data availability and research reproducibility in hydrology and water resources, *Scientific Data*, 6, 190030, <https://doi.org/10.1038/sdata.2019.30>, 2019.
- Strömqvist, J., Arheimer, B., Dahné, J., Donnelly, C., and Lindström, G.: Water and nutrient predictions in ungauged basins: set-up and evaluation of a model at the national scale, *Hydrolog. Sci. J.*, 57, 229–247, 2012.
- Tang, T., Stokal, M., van Vliet, M. T., Seuntjens, P., Burek, P., Kroeze, C., Langan, S., and Wada, Y.: Bridging global, basin

- and local-scale water quality modeling towards enhancing water quality management worldwide, *Curr. Opin. Env. Sust.*, 36, 39–48, 2019.
- Tilman, D., Fargione, J., Wolff, B., D’antonio, C., Dobson, A., Howarth, R., Schindler, D., Schlesinger, W. H., Simberloff, D., and Swackhamer, D.: Forecasting agriculturally driven global environmental change, *Science*, 292, 281–284, 2001.
- Toming, K., Kutser, T., Laas, A., Sepp, M., Paavel, B., and Nõges, T.: First experiences in mapping lake water quality parameters with Sentinel-2 MSI imagery, *Remote Sensing*, 8, 640, <https://doi.org/10.3390/rs8080640>, 2016.
- UN-Water: Summary Progress Update 2021: SDG 6 – water and sanitation for all, available at: <https://www.unwater.org/publications/summary-progress-update-2021-sdg-6-water-and-sanitation-for-all/>, last access: 7 September 2021.
- United Nations Environment Programme: GEMStat database of the Global Environment Monitoring System for Freshwater (GEMS/Water) Programme, International Centre for Water Resources and Global Change, Koblenz, available upon request from GEMS/Water Data Centre at: <https://gemstat.org/>, last access: 16 November 2020.
- United States Geological Survey: Water Quality Portal, available at: <https://www.waterqualitydata.us/portal/>, last access: 16 November 2020.
- Virro, H. and Kmoch, A.: GRQA code supplement, Zenodo [code], <https://doi.org/10.5281/zenodo.5082147>, 2021.
- Virro, H., Amatulli, G., Kmoch, A., Shen, L., and Uemaa, E.: Global River Water Quality Archive (GRQA), Zenodo [data set], <https://doi.org/10.5281/zenodo.5097436>, 2021.
- Wellen, C., Kamran-Disfani, A.-R., and Arhonditsis, G. B.: Evaluation of the current state of distributed watershed nutrient water quality modeling, *Environ. Sci. Technol.*, 49, 3278–3290, 2015.
- Welty, E., Zemp, M., Navarro, F., Huss, M., Fürst, J. J., Gärtner-Roer, I., Landmann, J., Machguth, H., Naegeli, K., Andreassen, L. M., Farinotti, D., Li, H., and GlaThiDa Contributors: Worldwide version-controlled database of glacier thickness observations, *Earth Syst. Sci. Data*, 12, 3039–3055, <https://doi.org/10.5194/essd-12-3039-2020>, 2020.
- Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., et al.: The FAIR Guiding Principles for scientific data management and stewardship, *Scientific Data*, 3, 160018, <https://doi.org/10.1038/sdata.2016.18>, 2016.
- Wood, E. F., Roundy, J. K., Troy, T. J., Van Beek, L. P. H., Bierkens, M. F., Blyth, E., de Roo, A., Döll, P., Ek, M., Famiglietti, J., Gochis, D., van de Giesen, N., Houser, P., Jaffé, P. R., Kollet, S., Lehner, B., Lettenmaier, D. P., Peters-Lidard, C., Sivapalan, M., Sheffield, J., Wade, A., and Whitehead, P.: Hyperresolution global land surface modeling: Meeting a grand challenge for monitoring Earth’s terrestrial water, *Water Resour. Res.*, 47, W05301, <https://doi.org/10.1029/2010WR010090>, 2011.
- Wu, Y. and Chen, J.: Investigating the effects of point source and nonpoint source pollution on the water quality of the East River (Dongjiang) in South China, *Ecol. Indic.*, 32, 294–304, 2013.
- Xu, X., Ester, M., Kriegel, H.-P., and Sander, J.: A distribution-based clustering algorithm for mining in large spatial databases, in: Proceedings 14th International Conference on Data Engineering, Orlando, FL, USA, 23–27 February 1998, IEEE, 324–331, 1998.