

# 1 Multiorder Hydrologic Position in Europe as a Set of 2 Metrics in Support of Groundwater Mapping at 3 Regional and National Scales

4 Maximilian Nölscher<sup>\*, 1</sup>, Michael Mutz<sup>2</sup>, and Stefan Broda<sup>1</sup>

5 <sup>1</sup>Federal Institute for Geosciences and Natural Resources (BGR), Sub-Department: Basic information Groundwater  
6 and Soil (B2.2), Berlin, 13593, Germany

7 <sup>2</sup>independent researcher

8 <sup>\*</sup>corresponding author: Maximilian Nölscher (max-n@posteo.de)

## 9 ABSTRACT

10 This dataset (EU-MOHP v013.0.1) provides information on the multiorder hydrologic position of a geographic point within its respective river network and catchment. More precisely, it comprises the three measures “lateral position” as a relative measure of the position between the stream and the catchment boundary/ watershed, “divide stream distance” as an absolute distance measure that serves as a proxy for the position within the catchment and “stream distance” as an absolute measure of the distance to the nearest stream. These three measures were calculated for several hydrologic (stream) orders and therefore reflect different spatial scales. Its spatial extent covers major parts of physiographical Europe and all of the 39 countries in European Economic Area (EEA39). Although there might be many potential use cases, this dataset serves predominantly as valuable static geophysical or environmental predictor variable among other input data for mapping or modeling tasks in the context of hydrogeology and hydrology.

## 11 1 Background & Summary

12 In recent years, data science tools such as machine learning are increasingly applied to and specifically developed for  
13 hydro(geo)logical challenges and research questions (!source: Zounemat 2020). In the field of hydrogeology, machine learning  
14 has been used successfully for groundwater level prediction and a variety of mapping tasks (!source: Wunsch 2021; Knoll,  
15 Desimone). Since machine learning models – except for hybrid- or physics-guided models – are purely based on data with  
16 no built-in knowledge of physical processes, it is important to provide as many features (synonyms: predictor variables,  
17 explanatory variables) as possible that have an impact on the target variable to potentially enable the machine learning algorithm  
18 to reproduce the result of the underlying process. For surface and near-surface processes, this criterion may be more or less  
19 satisfiable through the availability of remote sensing data, whereas for modeling subsurface processes such as in hydrogeology,  
20 this poses a serious challenge.

21 The key motivation of this dataset is to provide a set of features that introduce hydrological context to machine learning  
22 models regarding the horizontal position of a point within its catchment. Therefore, it functions as a proxy for multiple  
23 geophysical characteristics of a hydrologic system. It complements commonly available data sets and tackles the above  
24 mentioned challenge. This dataset is strongly inspired by [!source: Beelitz et. al.] and adapts their ideas and methods to  
25 the “EU-Hydro - River Network Database” but – in contrast – with purely free open source software and a strong focus  
26 on reproducibility. [!source: Beelitz et. al.] provides a comprehensive explanation of the motivation as well as a detailed  
27 discussion for further reading.

28 In their study, Beelitz et al. (2019) also provide the results from case studies to prove that the multiorder hydrologic position  
29 is a valuable feature when mapping diverse geophysical targets using machine learning. Its benefit to the performance of machine  
30 learning models has also been acknowledged by several other studies [!source: Using Boosted Regression Tree Models to  
31 Predict Salinity in Mississippi Embayment Aquifers, Central United States; Machine Learning Predictions of pH in the Glacial  
32 Aquifer System, Northern USA; Machine-learning models to map pH and redox conditions in groundwater in a layered aquifer  
33 system, Northern Atlantic Coastal Plain, eastern USA; The relation of geogenic contaminants to groundwater age, aquifer  
34 hydrologic position, water type, and redox conditions in Atlantic and Gulf Coastal Plain aquifers, eastern and south-central  
35 USA].

36 Being a static geophysical catchment attribute, the EU-MOHP data set can be used as features in any machine learning task  
37 in the domain of hydrology and hydrogeology. Because of the calculation based on the different stream orders, this data set can

be applied for multiple spatial scales and depths – from local via regional to continental scales. Examples of use cases might be the mapping of hydrogeochemical parameters or hydraulic variables like depth to groundwater, the prediction of groundwater levels or catchment classification tasks using unsupervised machine learning methods.

The EU-MOHP v013.0.1 data set comprises the 3 measures

- lateral position (lp)
- divide stream distance (dsd) and
- stream distance (sd)

for each of the 6 streamorders which leads to  $n_{measures}n_{streamorders} = 18$  different metrics to be used as features. Spatially, the data set covers major parts of physiographical Europe and all of the 39 countries in European Economic Area (EEA39). More precisely, it covers the 2 largest coherent land masses of the EEA39 (see fig. 2).

The generation of this data set is based on two basic data products which are the “EU-Hydro – River Network Database” and “EU-Hydro – Coastline” with the advantage that the dependencies are low from a data point of view. Therefore, it is possible to transfer the methodology to other regions with only little effort. For the calculation of the MOHP metrics in general, the following 4 essential data layers are required:

1. **river network** as vector layer with linestring geometries
2. **surface water bodies** such as lakes as vector layer with polygon geometries
3. respective **river basins** as vector layer with polygon geometries
4. **coastline** as vector layer with linestring geometries

These data layers are the basis for all further calculations of this data set. In this case, they were all obtained from the mentioned data sources provided by Copernicus Land Monitoring Service. The data set was manually downloaded from the website of Copernicus - Land Monitoring Service (“EU-Hydro - Coastline” 2021) and subsequently processed in the R programming language (see fig. 1; R Core Team 2020). Due to the memory size of this data set as well as for the sake of computational speed, a PostgreSQL database with PostGIS extension was used for some processing steps of vector data and GRASS GIS database was used for all final calculations of the MOHP metrics (see fig. 1D and E). For reproducibility and control reasons, all processing steps including the databases were tracked and executed through a data processing pipeline using the targets package in R (see fig 1; Landau 2021). More details can be found in [Methods](#) or in the code itself (see [Code availability](#)).

## Methods

The methods will be described by following the chronological order of the processing steps in the targets pipeline as it can be found in the code. This order is mainly guided by the dependencies of the processing steps. To better refer to the code, each processing step is provided with a paragraph naming its respective target names in the targets pipeline. The single processing steps are thematically summarised. The targets pipeline is split across multiple files in the `targets` directory (see `??`, line 28) and merged in `_targets` (see `??`, line 39). For more details, see these files. The targets marked (`helper target`) are not described here because they are not relevant to the methods and exist only for technical reasons.

### 1.1 Overview

The methods of this study conceptually follows the methods described in Belitz et al. (2019) which we highly recommend for further reading. The

@ref(fig:projectdirtree

### 1.2 Data Acquisition

The required river network data “EU-Hydro – River Network Database” as well as the coastline were manually downloaded from the Copernicus - Land Monitoring Service website (“EU-Hydro - Coastline” 2021, 2021). This step relates to fig. 1A. The unzipped files have a size of approximately 14 GB.

### 1.3 Processing Steps

All following steps relate to fig. 1C, D, and E.

#### 1.3.1 Data import

In this study, we used the coordinate reference system (CRS) ETRS89-extended / LAEA Europe with the EPSG code 3035. Therefore, all data was reprojected to this CRS after importing in case it differed.

## 85 River networks

### 86 Detailed Workflow

87 In the following, the description of the methods is oriented towards the structure of the targets pipeline to easily relate the  
88 methods description here to the source code in the repository. All steps required to understand the workflow will be described,  
89 for further details we refer to the source code.

#### 90 *Step 1: Data Acquisition*

91 The “EU-Hydro - River Network Database” was manually downloaded from <https://land.copernicus.eu> (for  
92 detailed link see references) as version v013. All downloaded and unzipped files have approximately 14 GB. The `!!river` is the  
93 only underlying data for the generation of the EU-MOHP dataset.

### 94 Hardware

95 The pipeline to generate the data set was executed on a DELL PowerEdge C4140 Server with an Intel Xeon Gold 6240R CPU  
96 and 384 GB installed RAM. The installed operation system is Microsoft Windows Server 2019 Standard, version 10.0.17763  
97 Build 17763.

98 [what is different to the Beelitz Paper and why] NHDPlusV2 data No pathlevel column Criterion to exclusively use free  
99 open source software

## 100 Data Records

101 Text.

## 102 Technical Validation

103 Text.

## 104 Usage Notes

105 Text.

## 106 Code availability

107 All processing and analysis was conducted with free open source software. All processing steps except for the download of  
108 the “EU-Hydro - River Network Database” (for practical reasons also referred to as river database) that was done manually  
109 are controlled and executed from within a targets pipeline in the programming language R [see fig. `!!source`]. Targets is an R  
110 package that provides a toolkit for reproducible workflows [`!!source`]. Spatial vector data such as the EU-Rivers are processed  
111 partly in R and a PostgreSQL (version 13) database with PostGIS (version 3.1.0) extension database for speed and memory  
112 reasons. For the same reason, all major raster calculations were conducted in a GRASS GIS (version 7.8.5-2) database. The  
113 database connections and all calculations in the databases are also controlled by this pipeline. For reaching a maximum of  
114 reproducibility, a docker container is provided to rerun all calculations easily. The R package `renv` is used for keeping track of  
115 the required R package versions and fits well to the combination with targets and docker to endure reproducibility.

## 116 Acknowledgements

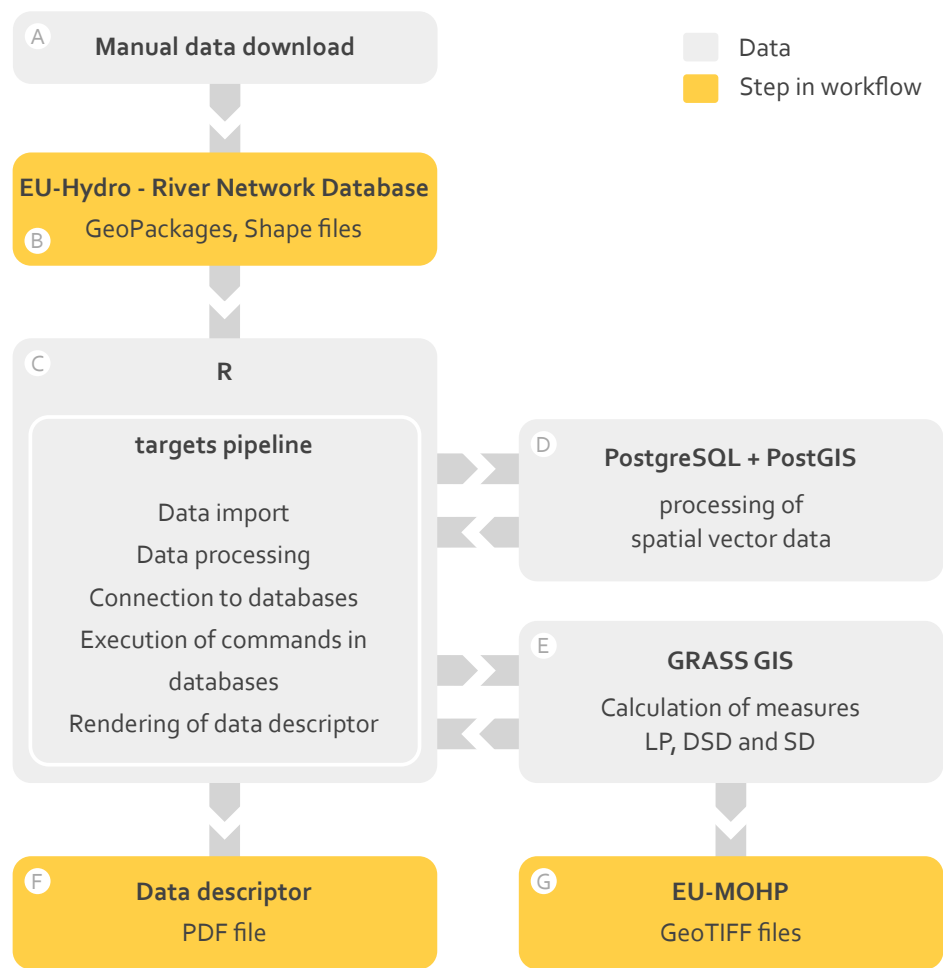
117 Text.

## 118 Author contributions statement

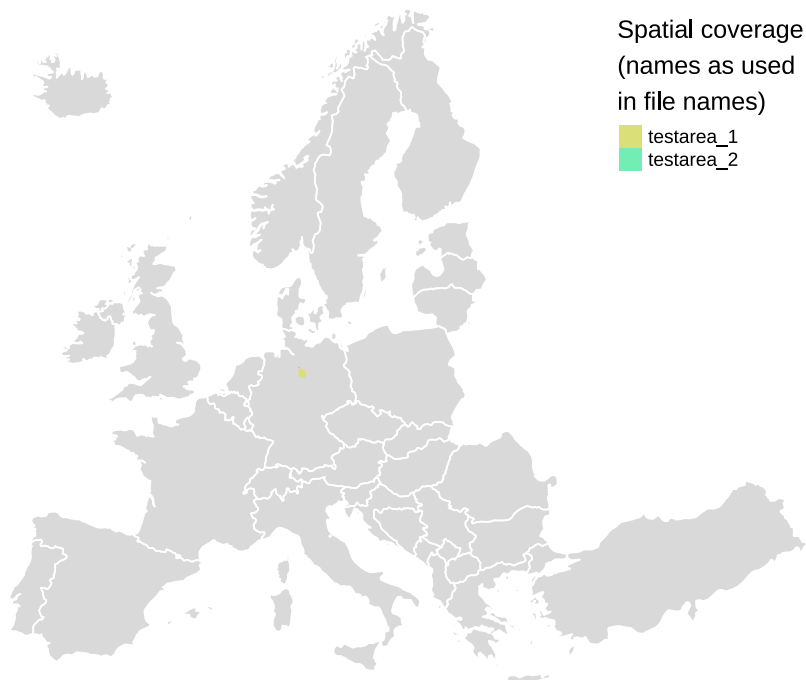
119 Text.

## 120 Competing interests

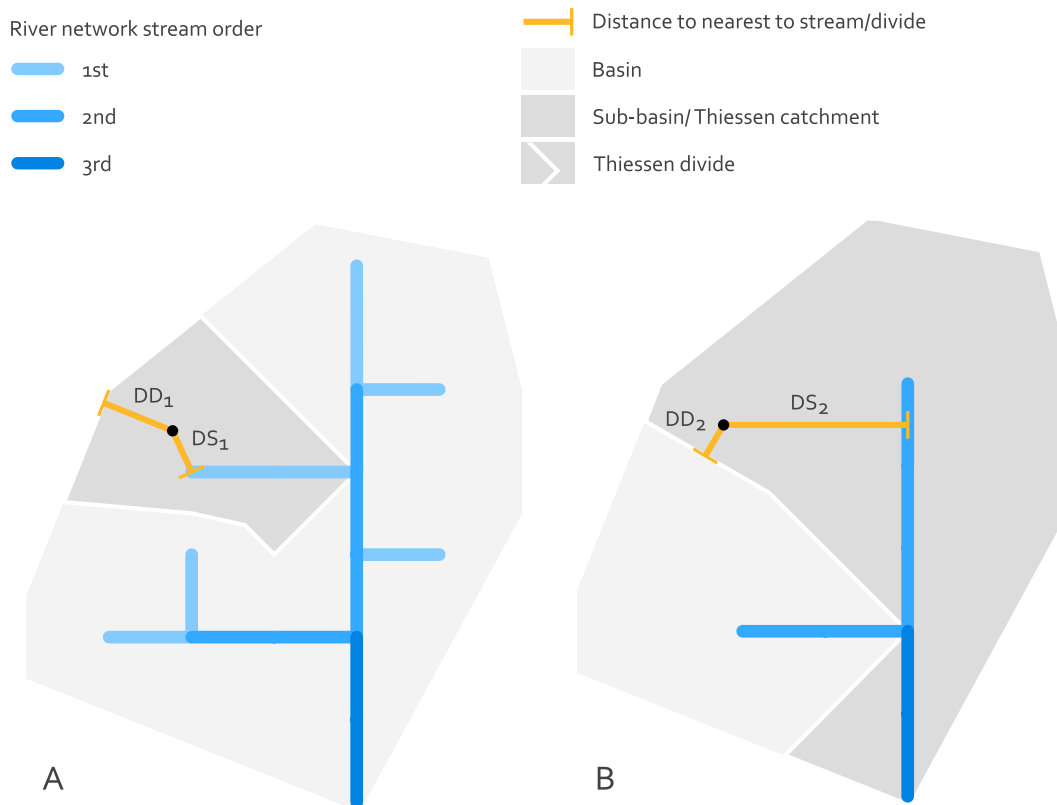
121 Text.



**Figure 1.** Workflow of the data processing in different software.



**Figure 2.** Spatial coverage of this data set. The legend labels show the names according to the file name. If you want to more precisely check whether your study area or area of interest is covered by this dataset, please visit [!!link to github readme](#).



**Figure 3.** Schematic representation of MOHP measures using two examples for the stream orders one (A) and two (B)

- 10 Belitz, Kenneth, Richard B. Moore, Terri L. Arnold, Jennifer B. Sharpe, and J. J. Starn. 2019. "Multiorder Hydrologic Position in the Conterminous United States: A Set of Metrics in Support of Groundwater Mapping at Regional and National Scales." *Water Resources Research* 55 (12): 11188–207. <https://doi.org/10.1029/2019WR025908>.
- "EU-Hydro - Coastline." 2021. *EU-Hydro - Coastline* " Copernicus Land Monitoring Service. <https://land.copernicus.eu/imagery-in-situ/eu-hydro/eu-hydro-coastline?tab=download>.
- Landau, William Michael. 2021. "The Targets R Package: A Dynamic Make-Like Function-Oriented Pipeline Toolkit for Reproducibility and High-Performance Computing." *Journal of Open Source Software* 6 (57): 2959. <https://doi.org/10.21105/joss.02959>.
- R Core Team. 2020. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.

## References

1. EU-Hydro - Coastline (2021). Last visited on 02/02/2021.
2. EU-Hydro - River Network Database " Copernicus Land Monitoring Service (2021). Last visited on 03/22/2021.
3. Stackelberg, P. E. *et al.* Machine Learning Predictions of pH in the Glacial Aquifer System, Northern USA. *Groundwater* **59**, 352–368, [10.1111/gwat.13063](https://doi.org/10.1111/gwat.13063) (2021). Last visited on 05/05/2021.
4. Knierim, K. J., Kingsbury, J. A., Haugh, C. J. & Ransom, K. M. Using Boosted Regression Tree Models to Predict Salinity in Mississippi Embayment Aquifers, Central United States. *JAWRA J. Am. Water Resour. Assoc.* **56**, 1010–1029, [10.1111/1752-1688.12879](https://doi.org/10.1111/1752-1688.12879) (2020). Last visited on 05/05/2021.
5. Degnan, J. R., Lindsey, B. D., Levitt, J. P. & Szabo, Z. The relation of geogenic contaminants to groundwater age, aquifer hydrologic position, water type, and redox conditions in Atlantic and Gulf Coastal Plain aquifers, eastern and south-central USA. *Sci. The Total. Environ.* **723**, 137835, <https://doi.org/10.1016/j.scitotenv.2020.137835> (2020).
6. Zounemat-Kermani, M. *et al.* Neurocomputing in surface water hydrology and hydraulics: A review of two decades retrospective, current status and future prospects. *J. Hydrol.* **588**, 125085, [10.1016/j.jhydrol.2020.125085](https://doi.org/10.1016/j.jhydrol.2020.125085) (2020). Last visited on 05/05/2021.
7. R Core Team. *R: A Language and Environment for Statistical Computing* (R Foundation for Statistical Computing, Vienna, Austria, 2020).
8. Landau, W. M. The targets R package: a dynamic Make-like function-oriented pipeline toolkit for reproducibility and high-performance computing. *J. Open Source Softw.* **6**, 2959 (2021).
9. Belitz, K., Moore, R. B., Arnold, T. L., Sharpe, J. B. & Starn, J. J. Multiorder Hydrologic Position in the Conterminous United States: A Set of Metrics in Support of Groundwater Mapping at Regional and National Scales. *Water Resour. Res.* **55**, 11188–11207, [10.1029/2019WR025908](https://doi.org/10.1029/2019WR025908) (2019). Last visited on 08/17/2020.
10. DeSimone, L. A., Pope, J. P. & Ransom, K. M. Machine-learning models to map pH and redox conditions in groundwater in a layered aquifer system, Northern Atlantic Coastal Plain, eastern USA. *J. Hydrol. Reg. Stud.* **30**, 100697, [10.1016/j.ejrh.2020.100697](https://doi.org/10.1016/j.ejrh.2020.100697) (2020). Last visited on 08/05/2020.