

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

Maximilian Nölscher^{*, 1}, Michael Mutz², and Stefan Broda¹

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

²independent researcher

*corresponding author: Maximilian Nölscher (max-n@posteo.de)

ABSTRACT

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. 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.

1 Background & Summary

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

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

In their study, (!source: Beelitz et. al.) also provide the results from case studies to prove that the multiorder hydrologic position is a valuable feature when mapping diverse geophysical targets using machine learning. Its benefit to the performance of machine learning models has also been acknowledged by several other studies (!source: Using Boosted Regression Tree Models to Predict Salinity in Mississippi Embayment Aquifers, Central United States; Machine Learning Predictions of pH in the Glacial Aquifer System, Northern USA; Machine-learning models to map pH and redox conditions in groundwater in a layered aquifer system, Northern Atlantic Coastal Plain, eastern USA; 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).

Being a static geophysical catchment attribute, the EU-MOHP data set can be used as features in any machine learning task in the domain of hydrology and hydrogeology. 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

39 unsupervised machine learning methods.

40 **Methods**

41 All processing and analysis was conducted with free open source software. All processing steps except for the data download
42 that was done manually are controlled and executed from within a targets pipeline in the programming language R [!!source].
43 Targets is an R package that provides a toolkit for reproducible workflows [!!source]. Spatial vector data such as the !!rivers are
44 processed partly in R and a PostgreSQL database (version 13) with a PostGIS (version 3.1.0) extension for speed and memory
45 reasons. For the same reason, all major raster calculations were conducted in a GRASS GIS database (version 7.8.5-2). The
46 database connections and all calculations in the databases are also controlled by the targets pipeline. For reaching a maximum
47 of reproducibility, a docker container is provided to rerun all calculations with little effort. The R package renv is used for
48 keeping track of the required R package versions and combines well with targets and docker to endure reproducibility.

49 **Detailed Workflow**

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

53 **Step 1: Data Acquisition**

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

57 **Hardware**

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

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

63 [1](#)

64 **Data Records**

65 Text.

66 **Technical Validation**

67 Text.

68 **Usage Notes**

69 Text.

70 **Code availability**

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

80 **Acknowledgements**

81 Text.

⁸² **Author contributions statement**

⁸³ Text.

⁸⁴ **Competing interests**

⁸⁵ Text.

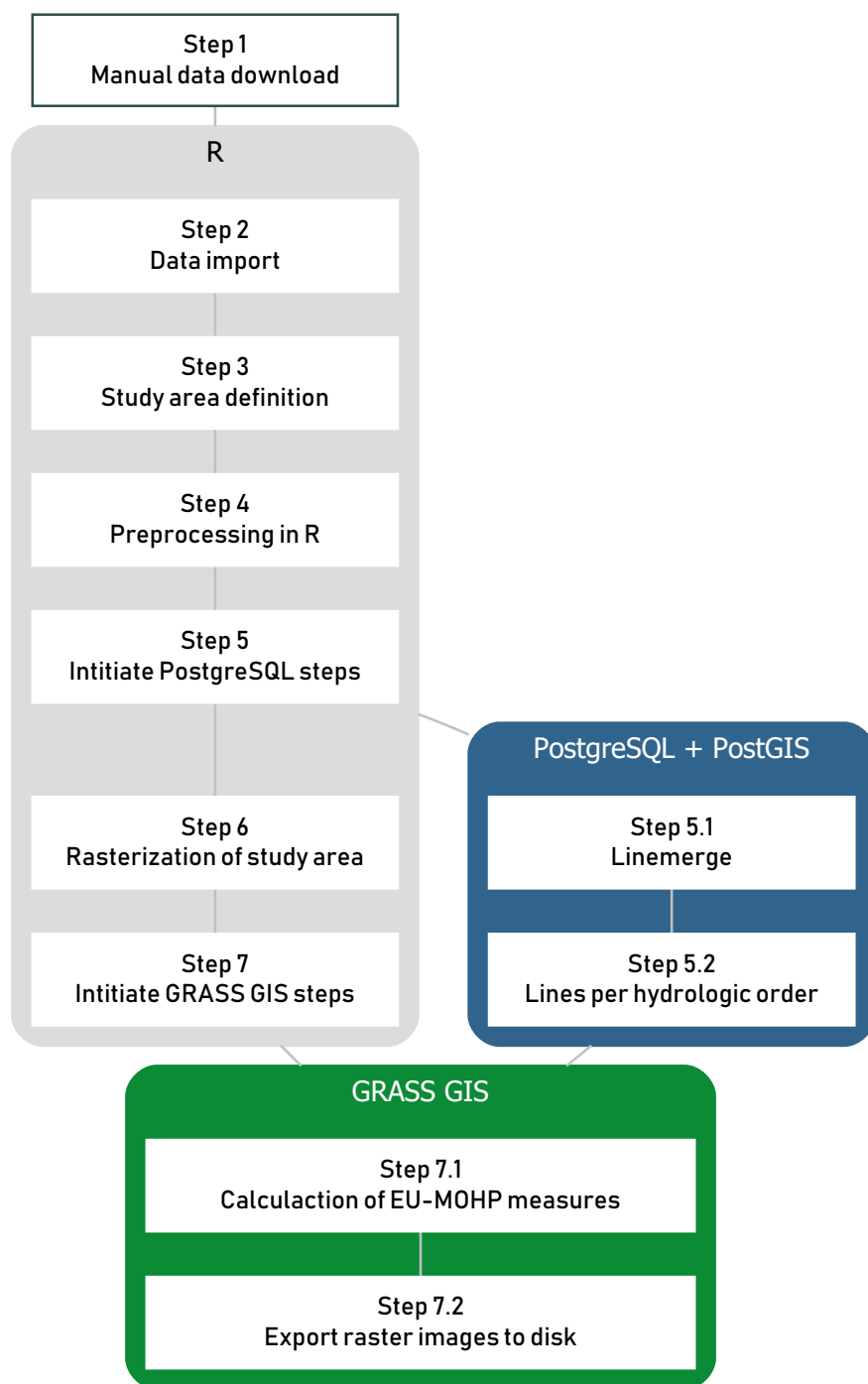


Figure 1. sdf

87 **References**

- 88 **1.** EU-Hydro - River Network Database â€™ Copernicus Land Monitoring Service (2021). Last visited on 03/22/2021.
- 89 **2.** R Core Team. *R: A Language and Environment for Statistical Computing* (R Foundation for Statistical Computing, Vienna,
- 90 Austria, 2020).

- 91 **3.** Landau, W. M. The targets R package: a dynamic Make-like function-oriented pipeline toolkit for reproducibility and
92 high-performance computing. *J. Open Source Softw.* **6**, 2959 (2021).
- 93 **4.** Belitz, K., Moore, R. B., Arnold, T. L., Sharpe, J. B. & Starn, J. J. Multiorder Hydrologic Position in the Conterminous
94 United States: A Set of Metrics in Support of Groundwater Mapping at Regional and National Scales. *Water Resour. Res.*
95 **55**, 11188–11207, [10.1029/2019WR025908](https://doi.org/10.1029/2019WR025908) (2019). Last visited on 08/17/2020.
- 96 **5.** DeSimone, L. A., Pope, J. P. & Ransom, K. M. Machine-learning models to map pH and redox conditions in groundwater in
97 a layered aquifer system, Northern Atlantic Coastal Plain, eastern USA. *J. Hydrol. Reg. Stud.* **30**, 100697, [10.1016/j.ejrh.](https://doi.org/10.1016/j.ejrh.2020.100697)
98 [2020.100697](https://doi.org/10.1016/j.ejrh.2020.100697) (2020). Last visited on 08/05/2020.