

## Review of

### Peckham et al. "HARBOR - Harmonized Attributes for River Basins in One Repo: Collated River Basin Data from Multiple Collections with a Software Toolkit"

([doi.org/10.5194/egusphere-2025-5786](https://doi.org/10.5194/egusphere-2025-5786))

#### General comments:

The authors collate a number of existing open data collections on river basin data from a range of public authorities in a single repository together with python functionalities to access the datasets. While the effort in bringing this data together and harmonizing the attributes is certainly commendable, the paper falls short of convincing me why it should be a research paper in HESS. If the main point is bringing the data together, I would rather see it as a contribution in a data journal such as ESSD with more focus on the harmonization efforts and the description of the differences. If it's about the research and model development that can be done with such a collection, I would expect at least one example of that as a main part of the research paper, to show how granting easy access to this dataset can advance our knowledge and understanding in hydrology. If it is about the river classification scheme I would expect a broader introduction and critical assessment rather than explanation of the existing schemes, a specific application and demonstration of the performance or benefits of the classification. And if it's meant to be about the python tools to access and harmonise the data, I'd rather expect a publication in GMD. As the paper is written now, it seems to be out of scope for HESS. I put some comments in the following nevertheless to elaborate some of the inconsistencies I experienced and would encourage the authors to consider these and maybe restructure/rewrite the paper in a way that fits better into one of the abovementioned options.

#### Specific comments:

##### 1) Section 2:

- What was the reasoning for including exactly these datasets? What were the criteria?
  - If it's US based, why is Caravan in there? If it's in there because of potentially differently calculated forcing data or processing, this should be stated clearly as the reason for including it.
  - Why CAMELS and not just using the CONUS catchments? If it's because of the meteorological data to calculate indices for hydrological similarity later, this should have been introduced as criteria from the start (and also why not simply existing meteorological datasets could be brought together with the larger data collections).
  - HYSETS: why is it included? Also, please cite the other datasets that you mention went into it.
  - MOPEX: here some papers that used the dataset are cited in the description. If this is an important feature, it should also be done for the other datasets. If it's not important, it should be left out as probably all datasets have been used by various authors. Also, the information that there were workshops to discuss things is non-relevant for the resulting dataset.
  - NOAA RFC:

- Is it relevant that this dataset includes forecasts? Was it included because of this? Then it should be explained. And what is HADS?
  - The specifics for the Missouri basin and the incomplete metadata subsection would fit better under a potential section where attribute harmonization/ curation of the whole data product would be described
- NSF CZO, LTER and NEON: why are these datasets in the collection? Because of the richness of other data that is being collected? Why is this important for the model application that was mentioned in the goals and the hydrological similarity classification?? Also, these entries contain a different level of detail in that individual columns etc. are described. This does not need to be part of a general overview.
  - USDA/ARS: STEWARDS database not title of the USDA reference. Where would I find the STEWARDS database then?
  - USDA SCAN/SNOTEL: why are the soil moisture/snow sites included?
  - USGS GAGES-II: again, level of detail in columns not fitting for the overview. But here, some criteria are listed. This would be a good general aspect for the whole collection.
  - USGS HLR: shouldn't this be listed in the second part of the paper? Also, it contains a very high level of methodological detail. If the methodology of defining Hydrologic Landscape Regions is the focus of this paper, this should be explained more, also in comparison to other strategies, in the main part about the region definition. Here it is a bit out of place.
- The datasets that went into the collection, after explaining the criteria for including them in the first place, could be organized in a more informative way than only with text blocks. For example, a table containing the variables that were selection criteria, temporal and spatial coverage, etc. Something that helps the reader to see the differences or overlap in the datasets and then the reasoning for combining exactly these datasets.
- 2) Section 3: Organization of HARBOR
- Please include the reference to the Repository at the beginning of that section where you describe what it contains.
  - The description itself reads rather as something that should be part of the README file on the repository than (in the main) part of a research paper.
  - Please also make clear how the individual datasets that go into the collection are licensed.
  - Where do the 53 Attributes come from, how where they defined and harmonized across the different datasets?
- 3) Section 4: Hydrologic similarity
- The purpose of this section is not clear. The defined terms seem very basic for the general audience of HESS. So why would they need to be explained at all? And if they need to be explained because they are integral part of the Hydrologic Landscapes that are supposed to be defined later, this should be made very clear.
  - Aridity index definition basically says there are many definitions, but not which one is supposedly used (if it is, because it's unclear from this paragraph) and why this was chosen.
- 4) Section 5: River Basin classification systems
- HLRs/SWB: Not clear why exactly these are chosen to be explained. Is the goal to find an especially good classification? Are these two the only ones?

- L. 336: which models? Is this an important criteria which models can deal with which classification?
- Was the extension of the SWB the goal from the beginning? Not checking which classification systems existed and evaluating the best for the model purpose? If SWB was the goal, a respective introduction, comparison and application of it should be considerably extended.
- Hydrograph-based classifications: why include them in the TSV file? What is the purpose of them? If Berghuijs already showed that they classify somewhat the same as the SWBs what is the additional information/benefit of including them?
- ML methods: this would be more of an application of the data collection? -> potentially in an application/outlook/further research section? As it is not used here and not compared in detail, it does not really fit into this section.
- Hydrologic Signatures: these need specific data to be calculated, right? how much of the HARBOR data collection contains these attributes? (a table here in the overview section would help to show this)

#### 5) Section 6

- Why did you apply the SWB to the GAGES-II CONUS? Reasons?
- You used a method to graphically close gaps in the classification? Does this also have a representation in the data/ a scientific reason? Or is there a valid reason for the gaps that is scientifically interesting?
- L. 415 f: file structure does not need to be in the main part of a research paper. Appendix or in README on the repository.
- L. 420. Why do you calculate the Budyko index? Just because it can be calculated from the existing data or because you use it for classification?
- Seasonality index: why do I need this? What's the purpose of including these calculations? Do I need it for classification? How?
- Paragraph l. 453 ff: this seems a bit random as a comparison. if you would want to check classification quality, wouldn't satellite imagery with identifiable vegetation be the way to go instead of specific forests? also "roughly match, "known distribution" (according to which source?), "appears to coincide" does not seem to be a serious validation effort
- l. 473: until here it has not become clear why I should care about/what's the scientific reason for this extended classification

#### 6) Section 7:

- L. 477: Is this now the goal of the whole paper? The requirements formulated here do not seem to be founded on the explanations in the sections before. If this is supposed to be the case, it should be linked better, such as (1) a test why between 12 and 20 classes are manageable for a model, (2) how considering hydrologic similarity improved your classification massively...
- Also there are quite some references to NextGen. Are these recommendations meant to be generally valid or only for NextGen application?

#### 7) Section 8:

- I'm not an extensive python user so I cannot comment on the various python functionalities. However, the description seems a bit excessive for a research paper -> maybe put in appendix?

#### 8) Section 9:

- This reads now that the main goal of the paper was improving the accessibility to the data via the python functions and the metadata harmonization. If the tools are the main point, this should potentially rather be a software publication in GMD or such. If

the data collection is the main point, a data description paper in ESSD might make more sense, with a focus on harmonization...

9) A1-A3: These paragraphs read as if they should rather be a part of the repository, along the specific datasets, and/or available at the original data source. I don't see the need for it in the paper.

A4: seems a bit obvious. Wouldn't it be sufficient to be clear on the respective coordinate system (if relevant, this could be part of the dataset overview table) and potentially offer conversions if necessary?

B: remove. Irrelevant here.

Technical corrections:

As I have many issues with the general scope, reasoning and structure of the paper, I did not go through every detail to check for typos. The issues that I noticed nevertheless are listed below:

I. 173: CHECK THIS - I guess this should have been checked before submission

I. 601: CHECK – same as above

I. 691: "500 feet" -> HESS submission guidelines: "metric system mandatory"