

Review of ‘LamaH | Large-Sample Data for Hydrology and Environmental Sciences for Central Europe’ by Klinger et al.
Submitted to Earth System Science Data.

Gemma Coxon

This paper describes the development of a large sample hydrology dataset for Central Europe (predominantly Austria). It compiles data across 859 gauges providing meteorological and hydrological timeseries as well as a vast array of catchment attributes covering climate, hydrology, land use, geology, vegetation and human influences. The timeseries are provided at both hourly and daily timescales.

The authors need to be congratulated on compiling a fantastic dataset that will be of great value to the hydrological community. It is a huge effort and it is very evident that a lot of care and thought has gone into the dataset. There is extensive discussion of the data sources and processing steps in the paper. There is good consideration of uncertainties and the figures are generally nicely presented.

I recommend publication of the dataset and paper, but the authors have some work to clarify the basin delineation and code availability. The paper is mostly well written but needs a really thorough proofread with quite a few sentences that are unclear (see list of technical corrections but this is not exhaustive!).

Main comments

Basin delineation and ‘headwater’ catchments. The delineation of the catchment boundaries is a key feature of the dataset but currently is not clear. I suggest the following:

- *Terminology* needs to be much clearer. What do you mean by ‘orographic catchment’ (is this a commonly used term? Do you mean topographic catchment?). I don’t believe you are using the term ‘headwater catchment’ correctly - I interpret headwater catchments as low order catchments found in the upper reaches of river basins. In which case statements like ‘In contrast, however, LamaH does not only consider headwater basins’ on L18 are not correct as CAMELS datasets also do not only consider headwater basins.
- *Data source* of the catchment boundaries – why did you use catchment boundaries from two products? It is not clear how the catchment boundaries are combined. L114 ‘As aggregation areas...’ This sentence doesn’t make sense and needs rewriting.
- *Figure 2.* I find this figure a little unclear and the figure caption is currently very long. It may be worth simplifying the figure and thinking about moving some of the explanatory text in the figure caption to the main text in the paper.

Potential evapotranspiration. The analysis in Section 4.2 is really interesting and an important addition to the paper. I understand the authors decision to not include the PET data of ERA-5,

but given the importance of this variable as forcing data for a large amount of hydrological models (particularly conceptual lumped hydrological models), it seems a shame not to include it as a variable. It is not entirely clear how you would derive PET from the reference evapotranspiration provided. Were other global PET products considered?

Code availability. Reproducibility for these large-sample datasets is key. The authors should consider making their code available alongside the dataset. I would also recommend a code availability statement to make clear the code that was used in the paper. For example, I believe you used Nans' code to calculate the hydrologic and climatic catchment attributes and it would be good to make this clear at the end of the paper.

Colour scales. Often diverging colour scales are used for sequential data (for example Figure 3a and 3b). I encourage the authors to change the colour scales on these plots to sequential colour scales as this is a more appropriate colour scale for sequential data values. There is a nice discussion of this issue in Section 3.2 of this preprint for HESS by Michael Stoelzle and Lina Stein: <https://hess.copernicus.org/preprints/hess-2021-118/>

Minor/Technical corrections

L20 and L57 'data basis' – unsure what is meant here. Can it be rewritten so it is clear?

L42 'are probably known to a broader audience' – can you make a more pertinent point here? What is significant about these particular missions?

L68 'and the United Kingdom' – this should be 'and Great Britain'

L89 'in nine different countries' – I would remove the word 'different' here.

Figure 1. Can you add the size of the circle to the legend in the plot making it clear how the circle size relates to catchment area?

Table 1. I don't think Table 1 adds much to the paper and would move to supplementary information.

L159 'respectively 61 runoff time series' – this doesn't make sense.

L165 It is not just changes in channel profile that lead to incorrect runoff calculation but also extrapolation of the rating curve, backwater effects etc. It may be worth expanding this a little and citing McMillan et al (2012) here (<https://doi.org/10.1002/hyp.9384>).

L171. When you aggregated the hourly timeseries to daily – what time period do you use? For example, daily flow timeseries in the UK is the mean river flow in a water-day, (09.00 to 08.59 GMT, for example; 09.00 1st December to 08.59 2nd December).

L184. Personally I would not interpolate any timeseries and leave this up to the data user.

Fig 3. Figure 3a legend title should be 'Start of continuous data record'. I also think the x-axis on the histogram in Figure 3b is incorrect – should go from 0 – 100%?

L190. I am unclear what you mean by 'gauge hierarchies'?

L270. 'for each of the 3 different basin delineations' – you don't need 'different' here and can be worded as 'for each of the 3 basin delineations'. Also L423, L601 and L602 – you don't need the word 'different' here.

L310. What do you mean by 'large notches in catchment shape'?

L555. What do you mean by 'herding'?

L605. I disagree that 'These uncertainties have been addressed'. They have been considered and discussed but I wouldn't say they have been addressed as many are not quantified in the dataset.