

Earth Syst. Sci. Data Discuss., referee comment RC2  
<https://doi.org/10.5194/essd-2022-38-RC2>, 2022  
© Author(s) 2022. This work is distributed under  
the Creative Commons Attribution 4.0 License.

## Comment on **essd-2022-38**

Patrick Bartlein (Referee)

---

Referee comment on "LegacyClimate 1.0: a dataset of pollen-based climate reconstructions from 2594 Northern Hemisphere sites covering the last 30 kyr and beyond" by Ulrike Herzschuh et al., Earth Syst. Sci. Data Discuss., <https://doi.org/10.5194/essd-2022-38-RC2>, 2022

---

General comments:

This paper describes a set of pollen-based climate reconstructions for the Northern Hemisphere from the LGM to present. The paper is obviously one of three, one describing the pollen data (Herzschuh et al., submitted, which I couldn't find), another describing the chronology (Li et al., 2022, ESSD-Disc), and this one, describing the reconstructions. There are obvious redundancies among the papers, and I think readers and potential users of the data will find it frustrating to have to track down three papers.

Overall, the paper is not that well organized, with motivations for some of the analyses (e.g. CCA) not appearing until the results section (Section 4, titled "Dataset assessment"), and tutorial material on the nature of pollen data as a palaeo-archive appearing in the discussion, as opposed to the introduction (and presumably also in the first paper of the series, which, with good cross-referencing among the papers, would make it superfluous here). Perhaps this disorganization arose in parting-out the papers.

There are several overarching issues and questions that should be addressed:

Why were January and annual temperature and annual precipitation chosen as the targets for reconstruction? A more appropriate set of climate variables might be those that mechanistically control vegetation like winter cold, summer warmth, and moisture stress. A lot of the paper is devoted to dealing with the obvious correlation between annual temperature and precipitation, but it is never actually established why this is an issue.

What was the role of the canonical correlation analysis? To simply explore the data

perhaps, but in fact it represents an alternative reconstruction approach. In any case, it's neither clear what the purpose of the analysis is, nor are the results fully explained.

The two reconstruction approaches, weighted-averaging – partial-least-squares (WA-PLS) and the modern analogue technique (MAT) may be frequently applied, but they are not without issues themselves. WA-PLS, as is the case with some other methods, tends to “compress” the reconstructions toward the center of the distribution of the climate data (see Liu et al, 2020, Proc. Royal Soc. A, <https://doi.org/10.1098/rspa.2020.0346>). This will reduce the amplitude of the time series of the reconstructions. MAT suffers from the no-analogue problem, typically diagnosed by looking at the dissimilarities. The performance of the two approaches are examined in Fig. 3, but there is no attempt to account for the obvious spatial patterns.

A number of the analysis steps are not explained much at all, with the results just briefly described before moving on. In particular, the significance testing in Section 4.2 isn't fully explained: What is the “take-home message”? What does this analysis say about the usefulness of the reconstructions.

The results are described in terms of mid-Holocene minus present (1.5 to 0.5 ka) long-term mean differences, and some unusual time series plots, but there is no attempt to assess the reasonableness of the reconstructions with respect to paleoclimatic first principles or to compare them with simulations or independent observations.

I think these issues are all basically addressable, and with a little overhauling (i.e. no new analysis, just more complete explanation and discussion), the paper(s) will make a useful contribution.

Specific comments:

line 62: “climate proxy synthesis studies”. Do you mean “syntheses of climate reconstructions” or “syntheses of climate proxies” (i.e. the pollen data)? It's the former that can be directly compared with climate-model output.

line 71: “The evaluation of climate model outputs...” It's actually the climate models that are being evaluated in data-model comparisons (of simulations and observations or reconstructions).

line 73: “strong changes in the climate driver” Are you alluding to changes in GHGs during the instrumental record? Changes in insolation, ice-sheet distribution and size, and GHGs between the LGM and present are much larger. For example, the companion CMIP

experiment to the LGM is the 4xCO<sub>2</sub> experiment. CO<sub>2</sub> has yet to double from pre-industrial levels yet.

line 74: "The extratropical Northern Hemisphere ... complex spatial and temporal ... patterns." Well, yes, but it's also where most of the pollen data is from. I don't think you need to motivate focusing on the Northern Hemisphere extratropics.

line 90: "Regarding the prevalence...". Just say "Pollen data from ... have been used..."

line 94: "high resolution". Temporal? Spatial? Also, the last millennium is part of the Holocene, and the late-Quaternary, so you might get some push-back from dendroclimatologists about this notion.

line 102: delete "the large" (I think we know extratropical Asia is large area.)

line 103: Whitmore et al. (2005) describes the modern pollen (and climate) data set for North America, not (paleo) precipitation reconstructions.

line 108: If "Herzschuh et al., submitted" is "LegacyPollen 1.0: A taxonomically harmonized global..." then how is that different from this paper (and the data sets on Zenodo)? Does it describe just the fossil-pollen data, or the modern data set too?

line 110: "Li et al., 2022). So there are three papers, 1) the pollen data set, 2) new chronologies, and 3) this paper, right? Why not just say that?

line 116: Why reconstruct temperature and precipitation, as opposed to climate variables that are mechanistically related to vegetation?

line 136: "For consistency with the amount (number?) of taxa...". This needs to be a little better explained. Why 70 taxa (except for tradition)?

line 147: "2000 km radius". Why 2000 km?

line 150: "metrics". Meaning something other than just the squared-chord distance?

line 151: "square-root transformed pollen percentages". It might be worth pointing out that the same transformation is embedded in the use of the squared-cord distance dissimilarity measure in the MAT approach.

line 156: "co-variation". Why is this an issue? It might be the case that covariation among predictands wouldn't be an issue if they were mechanistically related to vegetation, as in the case of variables like MTCO and GDD (Wei, et al., 2020, Ecology <http://dx.doi.org/10.17864/1947.194>)

line 161: "... partialling out the respective other variable". Please explain.

line 161: "We applied a Canonical Correlation Analysis...". What were the community, constraining, and conditioning matrices in this analysis? More to the point, what was the objective of this analysis?

line 164: "the ratio ... was determined...". Why and for what purpose?

line 191: Define "RMSEP" on first use in the text.

lines 190-220: What accounts for the spatial variations in RMSEPs? Data density? Data quality (of both the pollen and climate data)? Confounding environmental factors?

line 221: "significance test". Of what? What hypothesis does the Telford and Birks test address?

line 241: "we subtracted those means from every record". There are two mean values (6.5 to 5.5 ka and 1.5 to 0.5 ka), and "every record" implies to me the whole data set, LGM to present. Aren't you just looking at the difference between those two mean values? (And why 1.5 to 0.5 ka?)

line 243: "warmer and drier" Than what? (Which time period is the warmer and drier one?). Throughout this paragraph, the sense of change in climate has to be made explicit. For parallelism, you should adopt a standard way of expressing the changes, e.g. "warmer than present in the mid-Holocene" or "cooling from the mid-Holocene to present" but don't mix states and trend.

line 250: What's a "more gradual pattern"?

Figure 6: What exactly is plotted here? Why use a log age axis? An alternative depiction of all of the reconstructions, and their temporal and latitudinal variations would be a Hovmöller diagram.

Figure 8: I guess we're supposed to see that there are more correlation coefficients between temperature and precipitation close to zero in the "tailored" analyses. I've got nothing against violin plots, but I think a standard histogram would work a lot better.

line 301+: What are the implications of these statistics and their spatial patterns?

lines 315-343: This tutorial on pollen data, chronologies, etc. should probably be in the introduction, not the discussion.

line 378: "numerical mechanisms ... reduce the reliability" Please explain.

line 410: "TraCE 21k" is a transient experiment. The model used was CCM 3.

Code and data:

I was able to run the example R code without problems. However, the data sets, described and labelled (via the extension) as .csv files (comma-separated values), are instead tab-separated files, which usually have the extension ".tab", or sometimes ".txt". This situation prevents a user from getting a quick look at the data using a spreadsheet program.

P.J. Bartlein