

This is an ambitious and important study that, with revisions, should be published.

As a soil ecologist focusing on soil carbon–nutrient dynamics and their interactions under environmental and management changes, I found the work both impressive and valuable. Using comprehensive, high-resolution global datasets, the authors employ spatiotemporal machine learning approaches to map soil organic carbon (SOC), soil pH, soil type, and other variables at a remarkable 30 m resolution worldwide. They also estimate SOC changes over the past two decades. The harmonized legacy soil observations are particularly noteworthy, and the resulting maps will be valuable for diverse applications, such as driving soil carbon models, informing land management, and guiding policy decisions.

That said, there are several areas—both scientific and presentational—where improvements would substantially enhance the manuscript's rigor, clarity, and impact.

The manuscript uses an excessive number of abbreviations, which interrupts reading flow. Some abbreviations are not defined at first mention—for example, CCC and RMSE in the abstract. The authors should not assume all readers will be familiar with these terms. I recommend carefully reviewing the manuscript to ensure that all abbreviations are defined upon first appearance and that only essential abbreviations are retained. Given the already considerable length of the manuscript, the use of numerous abbreviations does not meaningfully reduce length and may hinder comprehension.

RE: Many thanks to the reviewer for his kind words. We have added explanations of abbreviations CCC, RMSE in the abstract now (also in the manuscript).

The manuscript is long enough to deter some potential readers. If journal policy allows, I recommend moving certain sections—such as extended details on data collection, preparation, modeling approach, and mapping techniques—into the Supplementary Information. This would allow the main text to focus more on:

Implications and potential applications of the dataset

Interpretation of key findings (e.g., variable importance, spatial patterns)

Broader impacts and future directions

While complexity is sometimes unavoidable, figures should be as simple, clear, and self-explanatory as possible. Specific suggestions:

Figure 1: Reorganize to emphasize the workflow; remove unnecessary logos. Use a consistent layout style (either top-down or left-right) with clearly separated blocks for each step.

Figure 2: The content can be succinctly described in a few sentences; consider removing the figure or replacing it with more impactful visuals.

Figure 5: Overly complex and difficult to interpret; despite repeated attempts, I could not fully understand it.

Figure 6: Contains too much information without adequate caption detail. Consider showing only the left panel with block diagrams; integrate textual explanations, interpretations, and distribution plots into the main text or Methods section.

RE: Regarding your request to move extended details on data collection, preparation, modeling approach, and mapping techniques into the Supplementary Information, we do not mind splitting and reducing paper and putting different parts into the supplementary materials. It appears, however, that the majority of ESSD papers do not use supplementary materials and there is no strict limit of the PDF size in terms of pages / total words (<https://www.earth-system-science-data.net/submission.html>). We have also contacted the editorial office of ESSD to get more guidance and help us decide (the question was whether to move different sections to supplementary material and they suggested “keeping all processing steps chronologically in the main text”). The journal editors suggested keeping all relevant descriptions in the main document. Our preference is also to keep all relevant technical information in the same document so that the readers do not have to jump between the main text and supplementary materials. Nevertheless, we have tried to make the paper easier to read and easier to locate and understand different sections. We fully agree with the reviewer that the first submission of our manuscript was long, with a lot of detail (and as such “long enough to deter some potential readers”). Production of this data, however, took over 2 years of dedicated work and in this period we discovered many new issues (some colleagues suggested that the paper could have been split into multiple articles). We hence find it important that the steps are described chronologically and without missing any important step that eventually influences results we get. We sincerely apologize to the reviewer for making such an extensive paper (60+ pages).

We have reworked Fig. 1 and added highlights of the key steps — it should be now clearer to readers. We prefer, however, to keep the logos as the data we use requires that we acknowledge the data sources / original data providing organizations.

Fig. 2 has been removed and we have instead inserted some text to explain steps.

Fig. 5 is indeed somewhat complex and unfortunately can not be simplified without as the steps are correct and relevant; because we are doing spacetime modeling, it is important to compare performance of models in time, spacetime and time only and hence the figure might seem difficult to follow, but is correct nevertheless.

Fig. 6 has been simplified following your recommendations.

Broader impacts and future directions are discussed in detail in subsection “Broader impacts and possible future development directions” on P52L1.

Scientific Comments:

Motivation for 30 m resolution. The rationale for mapping at 30 m resolution should be articulated more clearly. Higher resolution should serve a clear theoretical or practical purpose—for example, improving global SOC stock estimates, supporting fine-scale land management, or providing critical input to Earth system models. The current introduction touches on these but could better synthesize them into a concise, logically connected research

objective. The four research questions presented are somewhat disjointed; consider distilling them into a single, coherent framework.

RE: This is a good point. We have rewritten the introduction and added more clear motivation for using 30 m resolution on P4L16–34. The four research questions are analyzed and presented chronologically and we unequivocally answer each research question on P56L5 (“Conclusions”). Unfortunately, we do not understand the suggestion “consider distilling them into a single, coherent framework” — we consider our paper with its structure to be a single coherent framework, also illustrated in Fig. 1.

Multicollinearity of predictor variables. The manuscript does not address multicollinearity among predictors—a significant concern for high-dimensional datasets, especially with overlapping climatic and remote sensing variables (e.g., SAVI, NDVI, GPP, and climate metrics).

Multicollinearity can cause overfitting and obscure variable importance. The authors should clarify whether collinearity was assessed or controlled, and if not, explain the rationale.

Reducing redundancy could also decrease computation time and improve model interpretability.

RE: This is a relevant point and we have had a lot of discussion on this topic internally inside the group (indeed, theoretically speaking we could convert the long list of covariates to Principal Components and/or using sparse autoencoders to embeddings; see e.g.

<https://bradleyboehmke.github.io/HOML/autoencoders.html#sparse-autoencoders>; a problem with using embeddings or PCA, however, is that this adds another layer of modeling and this in fact increases significantly computing time as we would need to predict every pixel 2x). The model we use (Quantile Regression Random Forest) and corresponding framework with feature selection in scikit-learn (Repeated Subsampling-Based Cumulative Feature Importance — RSCFI; https://scikit-learn.org/stable/modules/feature_selection.html#recursive-feature-elimination), is in fact a robust framework and usually is over-fitting-proof, hence technically speaking — multicollinearity of variables we used is dealt with QRRF and RSCFI (“no part of the random forest model is harmed by highly collinear variables”); <https://stats.stackexchange.com/questions/168622/why-is-multicollinearity-not-checked-in-modern-statistics-machine-learning>). As we indicate on P17L19–24: “Feature selection would typically reduce the initial number of layers to 60–120, removing layers that marginally contributed to the final model” hence such combination of QRRF and RSCFI in addition helps reduce model complexity, without suffering from multicollinearity effects. In addition, by using original covariates, we are able to interpret the models and then detect which covariates we should put effort into maintaining in the future (e.g. post 2025). But it is a correct point — we start with a large number of covariates (300+) and many covariates overlap; using PCA or embeddings could help decrease multicollinearity and decrease model complexity, however this come at the cost of interpretability and an additional model (additional modeling step) is needed to fit e.g. sparse autoencoders. This actually increases computing significantly as we need to convert all pixels from original covariates to components (so in fact 2 rounds of predictions: 1st round to derive components/embeddings, 2nd round to generate predictions). We have added this discussion “OpenLandMap-soildb methods and data limitations” on P45L15 in the revised manuscript. If we are not able to deal with this issue in this paper, at least future work should consider the issue of multicollinearity / too many overlapping covariates.

Inclusion of bedrock depth. Bedrock depth is a critical factor influencing SOC stocks, particularly in mountainous regions where bedrock often occurs at shallow depths (<1 m). While the authors mention future inclusion, I suggest considering it now—global bedrock depth maps do exist and could be integrated relatively easily. Bedrock depth affects root distribution and carbon inputs to the soil, potentially altering model performance and the relative importance of predictors.

RE: Thank you for the comment. The lead author of the OpenLandMap-soildb paper is familiar with depth to bedrock i.e. have attempted mapping global distribution of the depth to bedrock (<https://doi.org/10.1002/2016MS000686>). Adding an accurate depth to the bedrock map at 30 m is high on our priority. The complexity of mapping depth to bedrock, however, is significant: we have only limited training (point data); depth to bedrock is a censored variable (<https://bookdown.org/mattdobra/Prelude/censoredcount.html#censored-models>) and hence we believe that modeling and mapping at 30 m resolution is not trivial (over-fitting, over-/under-prediction can often happen). To produce an accurate and detailed map of depth to bedrock could become a (multi-year) project in itself; we are unfortunately not able to produce a detailed map of depth to bedrock at the time-line of a few months. Previous maps of depth to bedrock we produced are only available at 1km or coarser spatial resolution, which might not be suitable to add for 30 m maps due to high heterogeneity. We have added this explanation in the discussion section on P44L3.

Overall, this is an excellent and timely study with strong potential impact. Addressing the concerns outlined above—particularly improving structure, clarifying motivations, simplifying figures, and addressing certain methodological points—will greatly strengthen the manuscript's readability and scientific contribution.

RE: We thank the reviewer one more time for his kind words and for suggestions to help improve the readability and clarity of the draft paper. To further enhance easy access to data, we have also prepared a simple app at: <https://world.soils.app>; we have also added a python tutorial on how to access data and derive prediction intervals is available from: https://github.com/openlandmap/soildb/blob/main/OpenLandMap_soildb_tutorial.ipynb. We hope that these types of interfaces and supplementary materials will help increase clarity of the methods and make the data easier to access and validate by anyone.