

Reply on RC2.

Dear editor and reviewers:

We would like to express our sincere gratitude to you for your thoughtful comments and constructive suggestions of our manuscript, which clearly help us improve the manuscript. Please find our replies below. The reviewer's comments are shown in black, and our responses are in red.

Kind regards,

Authors

Review of, “Mapping Antarctic Geothermal Heat Flow with Deep Neural Networks optimized by Particle Swarm Optimization Algorithm”, by Liu et al.

Review by Michael Wolovick

Summary

In this manuscript, Liu and coauthors infer Antarctic Geothermal Heat Flow (GHF) by training a neural network on a global database of heat flow measurements and a variety of geological or geophysical predictor variables. They train their model to predict spatially binned global observational GHF data using the input variables, without any knowledge of spatial structure (other than the spatial structure implicit in the predictor variables). They use a technique called particle swarm optimization to select the optimal hyperparameters in their deep learning model. They try four different configurations of this technique to ensure that their results are robust, and they also test the robustness of their model by leaving out variable amounts of input data in a densely sampled region of Europe. Their analysis reveals significant advantages of deep learning techniques as compared to linear regression. Their final predicted Antarctic GHF map is broadly reasonable from a geophysical perspective, with higher heat flow in West Antarctica and the Antarctic Peninsula compared with East Antarctica, although they do find some notable local maxima of GHF in East Antarctica under the Gamburtsev Subglacial Mountains and Vostok Subglacial Highlands. They conclude with a comparison of their model against four other published estimates, and against the limited in situ observational data available.

Major Comments

My expertise is in glaciology, numerical modeling, and geophysics, so I cannot truly evaluate the deep learning techniques that the authors have used. As far as I can tell, they appear to have been quite thorough from a machine learning perspective, with a lot of effort having been spent to ensure that their model fits the data well with the right set of hyperparameters. However, I do know something about Antarctic heat flow, and I do understand the importance of having good input datasets and predictor variables for a project like this, especially since their deep learning models have no internal knowledge of spatial relationships: their model treats each grid cell as an independent data point, and therefore any spatial structure in the output must come from spatial structure in the inputs. I also understand how vital it is that the input datasets contain realistic and accurate data underneath the Antarctic Ice Sheet; no matter how good the neural

network is, if it is supplied with bad inputs, then it will produce bad outputs. It is therefore vital that the authors manuscript shows readers maps of all of the predictor variables, both globally and with a south polar view for Antarctica.

Unfortunately, this manuscript does not show those maps. It is therefore impossible for readers to evaluate how much the authors' results can be trusted. These figures need not all be in the main text; with 16 variables (14 after removal of two highly colinear variables) there is a lot of information to display, so supplemental figures or an appendix would be fine here. However, it is impossible to properly evaluate the authors' model without seeing those inputs.

This is an especially important issue given the fact that at least two of the input variables are problematic in Antarctica. The "rock type" variable (which is one of only 3 that the authors show, in Fig 2), classifies the entire Antarctic Ice Sheet under the category of "ice". This of course makes sense for a dataset which represents the surficial rock type of a region, but we are interested in the heat flow underneath the ice. The subglacial rock type is generally unknown for most of Antarctica, so this dataset is worthless. Additionally, sediment thickness is another variable they use that is poorly constrained underneath most of the Antarctic Ice Sheet. While it is true that active-source seismic surveys have constrained sedimentary basin depths in a handful of locations, for the most part we simply have no idea how thick the sediment is underneath the ice, so this dataset is also worthless. The authors made decisions about which variables to keep and which to exclude on the basis of colinearity (sensibly choosing to discard redundant variables that were highly correlated with other variables), but they seemed not to have considered physical plausibility or under-ice uncertainty in their decision-making. Good input datasets for a project such as this one must be datasets for which the structure in Antarctica is well-constrained. For at least two of the inputs, that is categorically not the case. The authors need to look over their input datasets and remove those that are not well-constrained underneath the Antarctic Ice Sheet. Most obviously this includes rock type and sediment thickness, but they should double-check all of their inputs to ensure that they have realistic spatial structure in Antarctica. This will require that the authors retrain their deep learning models using only the datasets that are reliable in Antarctica. Unfortunately, this may degrade the quality of the fit and reduce predictive capacity in the rest of the world. However, that is simply the nature of the problem we are trying to solve. If the goal is to infer GHF in Antarctica, then there is no point in using datasets that are unconstrained there.

In addition to the requirement that the input datasets be well-constrained in Antarctica, it is also important that they be free of spatial artifacts, since the authors' deep learning model has no internal knowledge of spatial relationships; it treats every grid cell as an independent data point, and thus it relies on the input datasets to produce spatial structure. Unfortunately, the authors' output model (Fig 6) contains pronounced meridional stripes radiating out from the South Pole. This is likely a result of the fact that the authors interpolated all of their inputs onto a latitude-longitude grid with constant grid spacing. Constant grid spacing in lat-lon space works fine in the mid latitudes, but it can produce artifacts near the poles, and the authors' result clearly has such artifacts. Since the authors' deep learning model treats every grid cell as an independent data point, it follows that these meridional stripe artifacts in the output are a result of similar stripe artifacts in at least one of the inputs.

There are three main methods that they could use to fix this: 1) they could use projected x/y coordinates for their Antarctic prediction while keeping lat-lon coordinates for the rest of the globe, although this potentially introduces problems in applying a deep learning model trained on lat-lon data to a new set of x/y data if the statistical distributions of the two datasets are different; 2) they could use variable grid spacing in longitude, with more grid points in each row near the equator and fewer grid points in each row near the poles, a method that looks especially attractive given that their deep learning models treat the data as a list of independent points rather than a structured grid anyway; 3) they could keep their regular lat-lon grid but apply latitude-dependent smoothing in the longitude dimension in order to ensure that their input datasets have constant spatial resolution even as the grid converges near the poles¹. The exact method is up to the authors' choice, and they are of course free to choose a different method from the three that I propose here, but it is important that they appropriately pre-process their predictor variables to remove artifacts in polar regions, because their deep learning algorithm is not going to be capable of removing those artifacts on its own. And, of course, it is vital that the authors show us these predictor variables, so that we can verify for ourselves that they are indeed artifact-free.

My overall recommendation is that this paper needs major revisions. I chose major revisions rather than minor mostly because I am recommending that the authors retrain their models after removing datasets that are unconstrained in Antarctica and pre-processing to remove meridional artifacts. The manuscript itself might not need a great many

changes. The additional figures I requested showing the input datasets can be placed in a supplement or appendix rather than the main text, and most of the main text can probably be kept without too much change. It might very well be that the new model has a broadly similar distribution of GHF, just without the artifacts. However, I want to see the authors' models retrained after the changes to the inputs that I described above, and since I am recommending that the authors redo their main modeling work, I classify this as a major revision.

We sincerely thank the reviewer for the thorough and insightful evaluation of our manuscript. The reviewer's expertise in glaciology, numerical modeling, and geophysics has been invaluable in identifying critical issues that significantly improve the scientific rigor of this study. We have carefully addressed all major concerns as follows:

1. Display of Predictor Variable Maps. We fully agree that displaying all predictor variables is essential for readers to evaluate the reliability of our results. We have added Figure S1 in Supplementary Materials, which includes global maps of all retained predictor variables and their corresponding Antarctic polar stereographic views (EPSG:3031).
2. Removal of Unreliable Input Datasets. We agree with the reviewer's concern regarding the physical plausibility and uncertainty. We have removed rock type, sediment thickness and distance to hotspot from our predictor variables, as these datasets are indeed poorly constrained beneath the Antarctic Ice Sheet. The model has been retrained using only the reliable datasets. As anticipated, global predictive performance slightly decreased, but this trade-off is necessary and appropriate to ensure reliable Antarctic predictions.
3. Meridional Stripe Artifacts. We thank the reviewer for identifying this issue and for providing detailed suggestions for resolution. Upon investigation, we found that the meridional stripe artifacts were caused by an error in the projected coordinate system during visualization in ArcGIS. The original approach of directly interpolating in the EPSG:4326 coordinate system resulted in artificial patterns near the poles due to grid convergence. We have now corrected this by projecting the data from EPSG:4326 to EPSG:3031 (Antarctic Polar Stereographic Projection) before interpolation. The revised Figure 6 no longer exhibits these artifacts.

Minor Comments

L29-30: "As an important heat source beneath the Antarctic ice sheet, GHF directly affects the hydrological system under the ice sheet (Kang et al., 2022)."

While I appreciate the reference to a paper I am coauthor on, there is probably a better reference to use here. We didn't talk much about hydrology in that paper, although we did show basal melt rates.

Thanks for your advice. We have replaced the reference with a more appropriate citations(Siegert et al., 2016):

Siegert, M. J., Ross, N., Li, J., Schroeder, D. M., Rippin, D., Ashmore, D., Bingham, R. G., and Le Brocq, A. M.: Antarctic subglacial groundwater: a concept paper on its measurement and potential influence on ice flow, Geol. Soc. London Spec. Publ., 461, 197–213, <https://doi.org/10.1144/SP461.6>, 2018.

L33-34: "In addition, the complex interaction between GHF and climate results in a significant degree of variation in Antarctic ice mass distribution."

I'm not sure what exactly you mean here. How does GHF interact with climate? This sentence needs to be reworded or clarified.

L37-38: "...lays a significant factor for understanding the feedback mechanisms produced by Antarctic ice mass loss and predicting sea-level change"

This sentence also needs to be clarified.

L40-41: “However, the sparse and uneven distribution of in situ borehole data for GHF, coupled with the severe climatic challenges of direct measurements in the Antarctic continental interior, presents significant challenges for data acquisition (Fisher et al., 2015).”

145 This sentence should be rephrased. How does the sparse distribution of borehole data present a challenge to data acquisition? It would be more accurate to say that the challenges of data acquisition result in a sparse distribution of data. Perhaps rephrase as, “Unfortunately, the severe logistical difficulties involved in collecting direct measurements in the Antarctic continental interior ensure that the distribution of in situ borehole data for GHF is sparse and uneven
150 (Fischer et al., 2015).”

L42-48: “Conventional approaches fall into two categories: one based on the derivation of geothermal processes, such as decreasing west-to-east heat flow derived from some assumptions of geological conditions (Pollard et al., 2005), crustal and upper-mantle heat flow inferred from seismic models (Shapiro & Ritzwoller, 2004; Shen et al., 2020; Hazzard & Richard, 2024), and Curie temperature depths estimated using satellite magnetometry and thermal models (Maule et al., 2005; Martos et al., 2017). The other was from statistical methods such as multivariate similarity analysis (Stål et al., 2021), Bayesian inversion of multiple datasets (Lösing et al., 2020) and machine learning (Lösing & Ebbing, 2021).”

155 These sentences need to be reworked as well. It is wrong to describe the first set of sources as “deriv[ing] geothermal processes”. “Geothermal processes” is an ambiguous phrase that could be misinterpreted as referring to hydrothermal circulation, which none of these sources represent. In addition, many of the sources in the first category are also engaged in some form of statistical modeling, not process modeling. Shapiro and Ritzwoller, for example, use a similarity function to relate seismic structure in Antarctica to seismic structure elsewhere in the world, where GHF observations are available. They don’t perform any thermal modeling. It would be better to say that the first group use one type of data (usually seismic tomography or magnetic anomalies), which the second group use multiple types of
160 data. In addition, there are some references missing here.

165 Perhaps this section could be rephrased as: “Conventional approaches fall into two categories: on the one hand are those which use a single type of observation to infer GHF, most commonly seismic tomography (Shapiro & Ritzwoller, 2004; An et al., 2015; Lucazeau, 2019; Shen et al., 2020; Haeger et al., 2022; Hazzard & Richard, 2024) or magnetic anomalies (Maule et al., 2005; Purucker et al., 2012; Martos et al., 2017), although broad tectonic reconstructions have been used as well (Pollard et al., 2005). On the other hand, there are a newer set of statistical methods which integrate multiple types of observational constraints to infer GHF using multivariate similarity analysis (Stål et al., 2021), Bayesian inversion, (Lösing et al., 2020) or machine learning (Lösing & Ebbing, 2021).”

L60: “...deep learning algorithms ... due to its high accuracy...”
175 Should be “due to their high accuracy”.

(L29-60): We thank the reviewer for the detailed and constructive comments on the Introduction section. Based on the suggestions from both reviewers, we have substantially rewritten this section to address all identified issues:

180 “Geothermal heat flow (GHF) refers to the heat energy transferred from Earth’s interior to the surface via conduction or convection (Pollack et al., 2013). As a critical heat source beneath the Antarctic ice sheet, GHF not only directly affects the subglacial hydrological system and promotes basal melting, but also serves as an important boundary condition for numerical models predicting the Antarctic Ice Sheet(AIS) mass balance and global sea-level change (Obase et al., 2023; Pollard et al., 2005; Wearing et al., 2024; Llubes et al., 2006). Furthermore, characterizing the spatial distribution of GHF over Antarctica is crucial for comprehending the continent’s past and present tectonic evolution (Artemieva, 2011; Reading et al., 2022).
185

Unfortunately, severe logistical challenges associated with collecting direct measurements in the Antarctic interior have resulted in a sparse and uneven distribution of in situ borehole GHF data (Fisher et al., 2015). Conventional approaches fall

190 into two categories: on the one hand are those which use a single type of observation to infer GHF, most commonly seismic
tomography (Shapiro & Ritzwoller, 2004; An et al., 2015; Lucazeau, 2019; Shen et al., 2020; Haeger et al., 2022; Hazzard &
Richard, 2024) or magnetic anomalies (Maule et al., 2005; Purucker et al., 2012; Martos et al., 2017), although broad tectonic
reconstructions have been used as well (Pollard et al., 2005). On the other hand, there are a newer set of statistical methods
195 which integrate multiple types of observational constraints to infer GHF using multivariate similarity analysis (Stål et al.,
2021), Bayesian inversion, (Lösing et al., 2020) or machine learning (Lösing & Ebbing, 2021). While these approaches exhibit
consistency at the continental scales—characterized by higher GHF beneath West Antarctica and lower values in East
Antarctica—substantial discrepancies persist at regional scales. Methods relying on single observation types are typically
constrained by limited data resolution and spatial coverage, as well as by underlying assumptions that may lack universal
200 validity. For instance, seismic tomography-based approaches provide regional-average GHF estimates derived from data
with limited sensitivity to upper crustal composition and a coarse lateral resolution of 600–1000 km across Antarctica
(Shapiro & Ritzwoller, 2004; O'Donnell et al., 2019). As demonstrated by Goutorbe et al. (2011) and Lucazeau (2019),
integrating multiple observables yields more robust results than those derived from any single dataset. Specifically, Stål et al.
(2021) showed that using 14–19 sets of observables produces a misfit of less than 10 mW m^{-2} , whereas additional datasets
may introduce excessive noise without significantly improving estimates. Consequently, multi-observable approaches
205 necessitate a careful selection of features with adequate Antarctic coverage and strict control over the number of inputs.
Uncertainties in the original input data can propagate through the modeling process, and the resulting uncertainties in
subglacial GHF estimates can substantially impact ice sheet mass balance simulations. Given that Antarctic ice sheet dynamics
remain the largest source of uncertainty in future sea-level rise projections—with estimates for the year 2100 ranging from
–5 to 43 cm of sea level equivalent under high emission scenarios (Seroussi et al., 2020; IPCC, 2021)—reducing GHF
uncertainty is critical for improving the reliability of sea-level change predictions.

210 Recently, deep neural networks (DNNs) have emerged as powerful tools for synthesizing high-dimensional geoscience data,
leveraging their formidable nonlinear mapping capabilities. Their efficacy has been proven in improving estimates of Antarctic
ice sheet surface melt (Hu et al., 2021), predicting seasonal sea ice extent (Andersson et al., 2021), and simulating basal melt
rates beneath ice shelves (Burgard et al., 2022). However, current neural network models encounter two primary challenges.
215 First, the performance of DNNs is highly sensitive to numerous hyperparameters; manual or suboptimal tuning often leads to
poor generalization or overfitting. Second, as inherently opaque "black-box" models, DNNs seldom provide reliable
probabilistic estimates or confidence intervals. This lack of quantifiable uncertainty limits their applicability in downstream
earth system modeling where error propagation is a concern.

220 To address these issues, this study proposes a hybrid framework that couples DNNs with Particle Swarm Optimization (PSO)
algorithms to refine parameter selection, underpinned by a Bayesian module for robust uncertainty quantification. This
integrated approach introduces two key processes aimed at enhancing model generalization and reliability. First, the global
search capability of PSO is leveraged to optimize DNN hyperparameters, thereby minimizing the objective function and
improving predictive accuracy in data-sparse regions.. Second, the integration of a Bayesian module facilitates the
225 decomposition of uncertainty into aleatoric components (stemming from input data noise) and epistemic components (inherent
in the model architecture and parameters). In the following sections, we detail the dataset construction and methodology,
provide an analysis of discrepancies between the new GHF estimates and prior predictions, and discuss potential uncertainties
along with their implications for future investigations."

230 L91-92: " Subsequently, these filtered, high-quality point measurements were aggregated by calculating the mean value
within a $0.5^\circ \times 0.5^\circ$ latitude-longitude grid. "

See my major comments about the problems with using a regular lat-lon grid when studying the polar regions.

235 We thank the reviewer for raising this concern, which relates to the major comment about potential artifacts from using a
regular lat-lon grid in polar regions. As described in our response to the major comments, we have addressed this issue by
projecting the data from EPSG:4326 to EPSG:3031 (Antarctic Polar Stereographic Projection) before interpolation.

Additionally, we have visualized all input variables in Supplementary Figure S1, which confirms that the input datasets are free of meridional stripe artifacts in Antarctica.

Figure 1

Would it be good to include a couple sentences talking about the overall geographic distribution of the global data used to constrain the model? By eye, these data seem to be heavily biased towards wealthy countries, with much lower data density in Africa, South America, and the Middle East.

In addition, the color scale should be changed. Blue-white-red is appropriate for data that represent anomalies with respect to a mean or zero value. The GHF measurements being shown here are all positive, however, so a different color scale should be used.

We thank the reviewer for the suggestions regarding Figure 1. We have added a brief discussion of the geographic data distribution in the text:

“The filtered dataset exhibits significant spatial clustering (Fig. 1). Data coverage is substantially denser in North America, Europe, and parts of East Asia, corresponding to regions with longer histories of geothermal exploration. In contrast, Africa, South America, the Middle East, and Antarctica have markedly sparser coverage, with large areas containing few or no measurements. This geographic bias poses a challenge for empirical GHF modeling, as the training data are not an unbiased representation of Earth's geological diversity, and certain tectonic settings are overrepresented relative to others (Stål et al., 2022).”

Additionally, we have revised Figure 1 to use a sequential color scale that better represents the continuous positive range of GHF measurements.

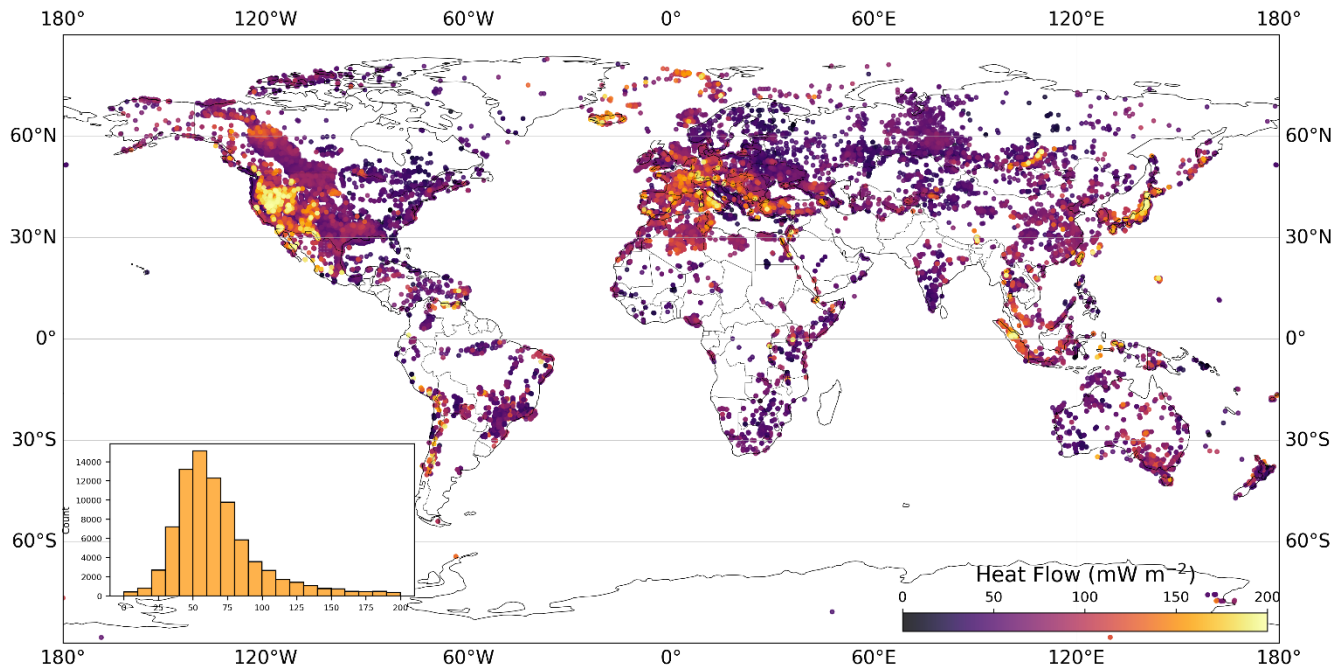


Figure 1. Spatial distribution of global GHF measurements used for model training.

Table 1

As I discussed in the major comments, all of these geophysical features need to be shown to the reader, both in global view and in south polar view. These figures can be placed in a supplement or appendix if necessary.

In addition, some of these data inputs have two sources listed. What does it mean when two sources are listed? Does that mean that the dataset is the mean of both sources? Or is it the case that one source is a publication and the other is a link to the actual dataset?

We thank the reviewer for these comments. As described in our response to Major Comment 1, we have added Figure S1 in the Supplementary Materials showing all retained predictor variables in both global view and Antarctic polar stereographic view (EPSG:3031). Some features list two data sources: the first reference refers to data from around the world, while the second provides Antarctic-specific data as a regional supplement. We have clarified this in the revised table:

" Note: Where two references are listed, the first provides global coverage and the second supplements with Antarctic-specific data offering higher regional resolution."

L115-116: "Sedimentary layers, due to their low thermal conductivity, act as an insulating blanket, significantly influencing the dissipation of deep-seated heat"

That may be true, but unfortunately, we have no meaningful constraint on sediment thickness underneath the ice sheet, at least not on a large scale. That is the challenge for a project like this: useful datasets are not merely those that have a meaningful physical relationship with heat flow, but those that have a meaningful relationship with heat flow and which are well-constrained in Antarctica. Excluding sedimentary thickness will, no doubt, reduce the quality of the global fit. However, the challenge of a project like this is to generate a model that can explain global heat flow using only variables that are known and well-constrained in Antarctica. Any predictive power added by sedimentary thickness will be of no help in Antarctica.

L123-125: "The Global Lithological Map (GLiM) database (Hartmann & Moosdorf, 2012) provides surface rock type data, explaining spatial variations in thermal conductivity."

Same concern as above. Their map (at least as shown in your Fig 2) lists the entire Antarctic Ice Sheet as the "ice" rock type, which is useless for inferring subglacial heat flow.

L115-116 and L123-125: We thank the reviewer for these important comments. We agree that the datasets for this project must not only have a meaningful physical relationship with heat flow, but also be well-constrained in Antarctica. As shown in Figure 2 (Rock Type), the entire Antarctic continent is classified as "Ice and Glaciers (IG)" in the Global Lithological Map database, which provides no information about the subglacial geology. Similarly, sediment thickness is poorly constrained beneath the Antarctic Ice Sheet on a continental scale. We have removed both rock type and sediment thickness from our predictor variables and retrained the model using only datasets that are well-constrained in Antarctica. As anticipated, this reduced the quality of the global fit slightly, but this trade-off is necessary and appropriate since any predictive power added by these variables would be of no help in Antarctica where they are unconstrained.

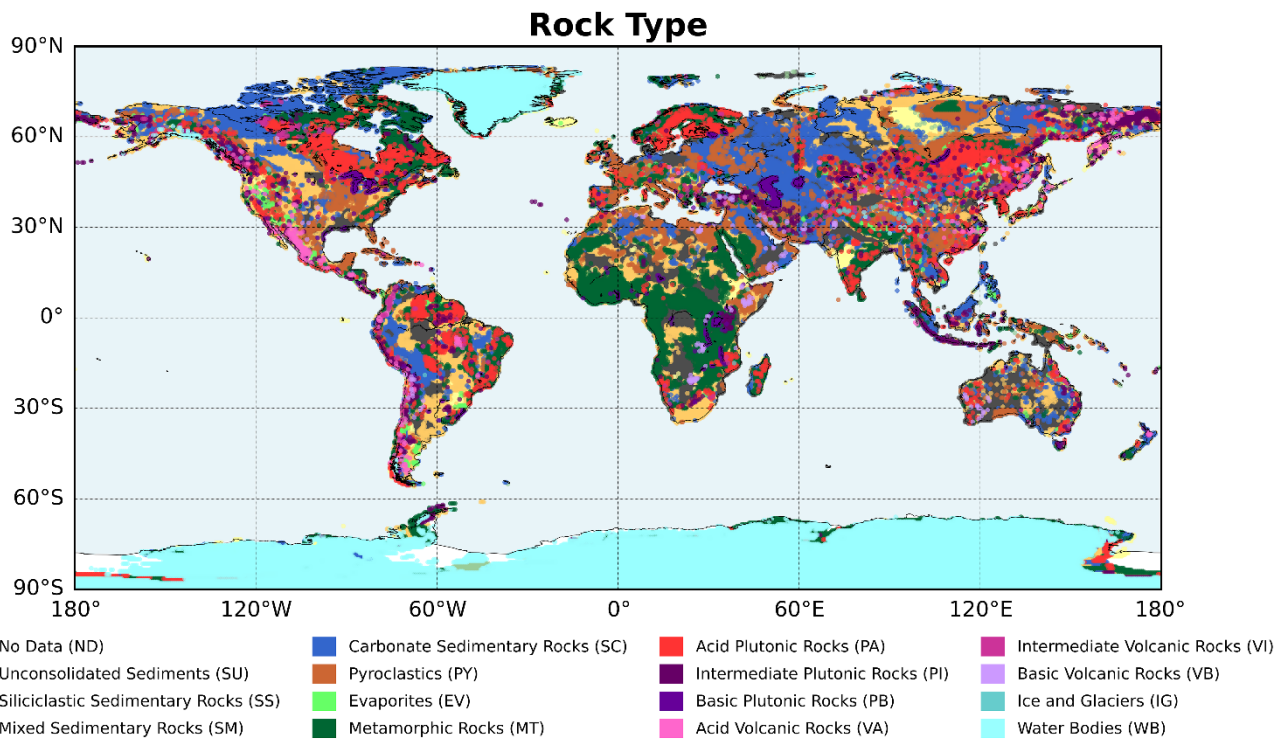


Figure 2. Global Lithological Map

L127-128: “To ensure dataset consistency, all predictor variables were resampled to a uniform $0.5^\circ \times 0.5^\circ$ grid using Ordinary Kriging.”

As I discussed in my major comment, a uniform lat/lon grid can produce meridional stripe artifacts near the poles. Potential solutions include: 1) using projected x/y coordinates in Antarctica; 2) using uneven grid spacing in longitude; 3) using latitude-dependent smoothing in the longitude dimension. Or perhaps a different solution that I haven’t thought of. But regardless, something has to be done to help this uniform lat-lon grid perform better near the South Pole.

We thank the reviewer for emphasizing this issue. As explained in our response to the Major Comments, we have addressed this problem by adopting the first solution suggested by the reviewer: using projected x/y coordinates for Antarctica. Specifically, we project the data from EPSG:4326 to EPSG:3031 (Antarctic Polar Stereographic Projection) before interpolation. The revised Figure 3 confirms that the Antarctic GHF predictions are now free of these artifacts.

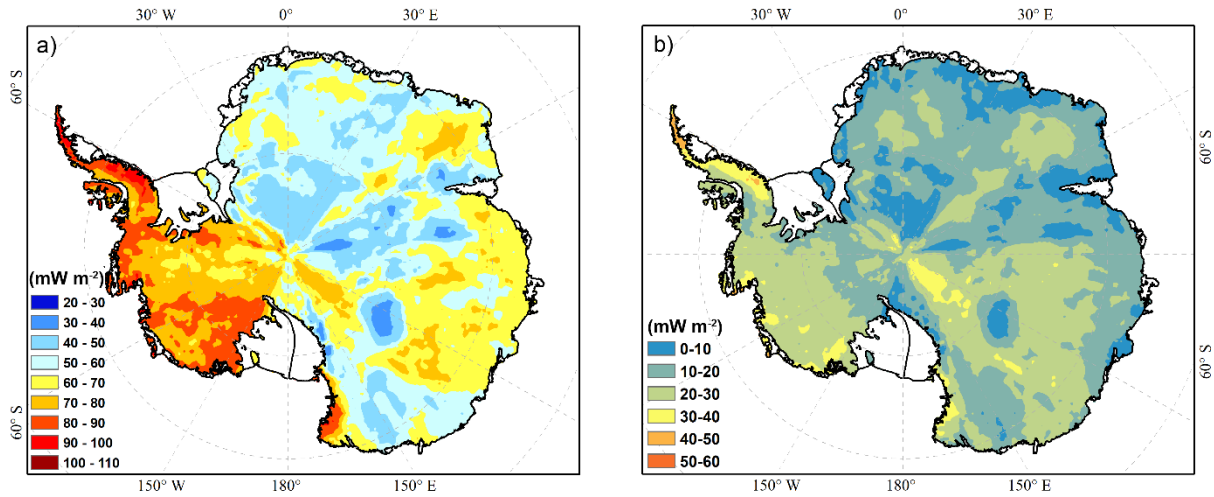


Figure 3. Predicted GHF distribution and associated uncertainty across Antarctica.

L157: “the Adam optimizer”

Does this need a reference?

We agree. We have added the appropriate reference for the Adam optimizer (Kingma & Ba, 2015).

Equation 3 I thought that R2 was the squared correlation coefficient? The formula for that would be:

$$R^2 = \left(\frac{\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2 \sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2}} \right)^2$$

Am I wrong about that? Is this a different definition of R2?

Indeed, there are two common definitions of R² in the literature:

(1) Squared Pearson Correlation Coefficient (as the reviewer described): $R^2 = r^2 = \left(\frac{\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2 \sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2}} \right)^2$

(2) Coefficient of Determination (used in our manuscript): $R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{SS_{res}}{SS_{tot}}$

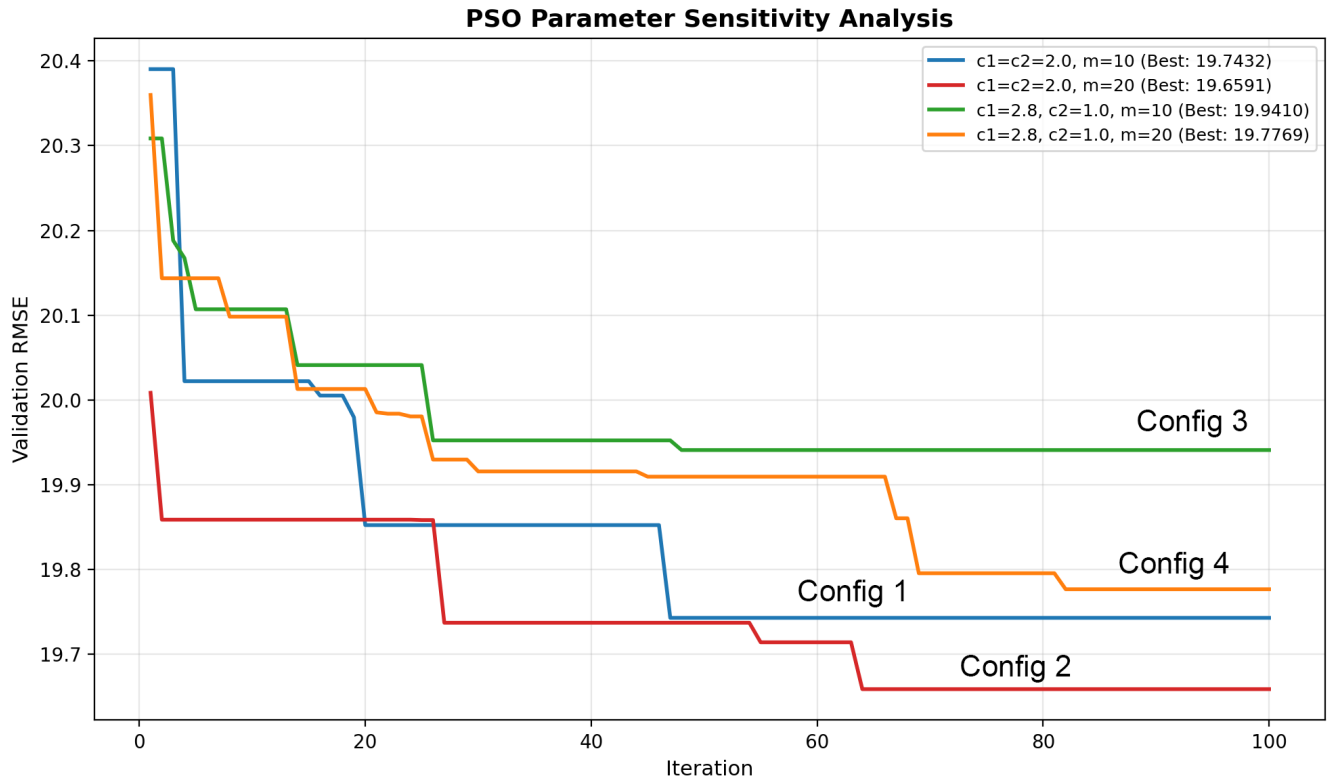
In our study, we adopted the second definition (coefficient of determination), which is the standard metric used in regression analysis and machine learning applications for evaluating model performance. Unlike the squared correlation coefficient, the coefficient of determination can yield negative values when model predictions perform worse than simply using the mean of the observed values as a predictor. A negative R² thus indicates that the model fails to capture the underlying patterns in the data.

We have revised the manuscript to explicitly clarify which definition of R^2 we employed and added an explanation that R^2 can become negative when the model predictions are worse than the mean baseline. Please refer to the revised Section 3.5 for details.

Figure 4

There is not much range on the y-axis here. Does that mean that all four of the configurations tested here have roughly the same performance? Or that the final result is relatively insensitive to the hyperparameters? In any case, the text should probably discuss the narrow range at some point.

You are right. We have added this discussion to the revised manuscript:



“Figure 4 illustrates RMSE trends across 100 iterations for the four configurations. Config2 ($c_1 = c_2 = 2.0, m = 20$) achieved the lowest final RMSE of 19.65 mW m^{-2} , followed by Config1 (19.74 mW m^{-2}), Config4 (19.77 mW m^{-2}), and Config3 (19.94 mW m^{-2}). The symmetric acceleration coefficients setting ($c_1 = c_2 = 2.0$) consistently outperformed the asymmetric configuration. Notably, the narrow performance differences (less than 0.3 mW m^{-2} between best and worst configurations) suggest that for this specific GHF prediction problem, the optimization landscape is relatively smooth,

allowing all configurations to converge to similar solutions. Nevertheless, Config2 was selected as the optimal configuration for subsequent model training based on its lowest validation RMSE."

Figure 5

Why does the circle enclosing your test region include parts of the Black and Aegian Seas? You have excluded marine observations from your dataset, so it seems like you could make a better dense test region by shifting the circle to only cover terrestrial parts of Europe.

In addition, why is R^2 negative for the linear regression model? Is this a function of the fact that you have defined R^2 differently than normal.

We have re-excluded data with the "marine" domain attribute from the IHFC database. For the NGHF database, we only retained data with geography codes A, B, C, D, E, F, G, and H (representing continental regions: Africa, North America, South America, Australia, Europe/Greenland, miscellaneous lands, Antarctica, and Asia/Arabia/India, respectively), excluding all oceanic measurements. However, some data points near coastlines may still appear because the original database classification includes "continental (lake, river, etc.)". Figure 1 has been updated to reflect this correction.

Regarding the test region: We have algorithmically re-adjusted the position of the test region circle to ensure it is more concentrated on areas with the highest density of terrestrial observations, minimizing overlap with marine regions.

And the negative R^2 value for the linear regression model is a direct consequence of using the coefficient of determination definition (as discussed in our response to the reviewer's comment on Equation 3). R^2 can become negative when the model's predictions perform worse than simply using the mean of the observed values as a constant predictor.

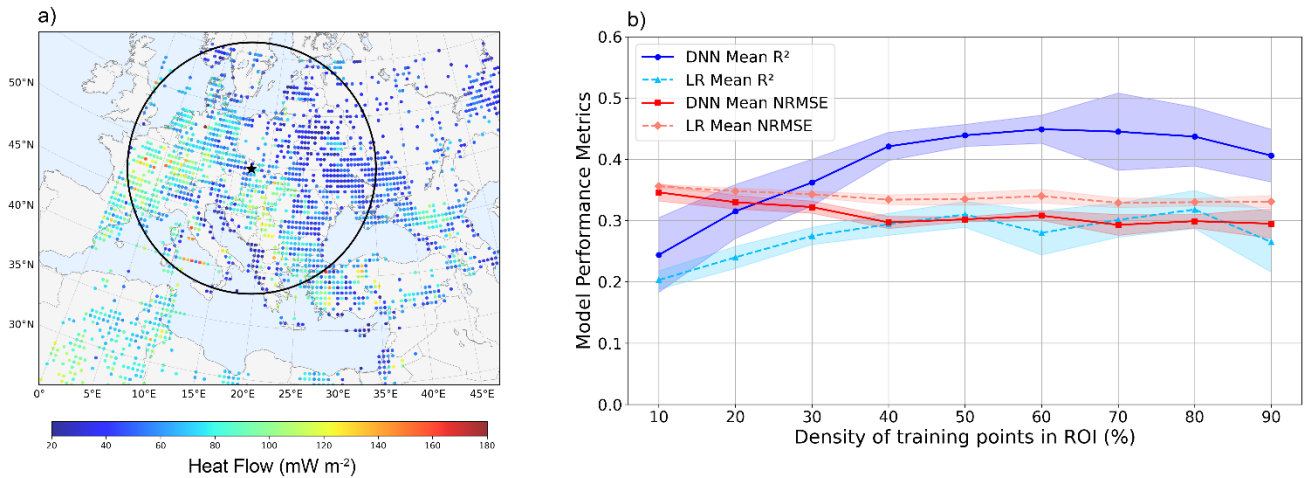


Figure 5. Performance of DNN and linear regression methods in experiments with different densities of ROI regions.

L318-322: “However, in the Gamburtsev Subglacial Mountains, Vostok Subglacial Highlands, and the area around Subglacial Lake Vostok, there is an increasing trend of heat flow values, which shows that these regions may have been affected by deep tectonic activity or localized heat sources (Artemieva, 2022).”

385 My own inversion for GHF (Wolovick et al, 2021) also showed a local maximum of GHF in the Gamburtsev Mountains which is necessary to fit observations of subglacial water networks there.

390 We thank the reviewer for this comment and for sharing the relevant findings from Wolovick et al. (2021). Following the reviewers' recommendations to remove predictor variables that are poorly constrained in Antarctica, we have retrained our model with the refined variable set. Unfortunately, the updated predictions show elevated GHF values in the area around Subglacial Lake Vostok, but the local maximum in the Gamburtsev Subglacial Mountains identified in the original submission is no longer apparent in the revised model. This difference may reflect the trade-off between global predictive power and Antarctic reliability: by excluding variables that are unconstrained beneath the ice sheet, we may have reduced the model's ability to resolve certain local features.

395 Figure 6

400 The meridional stripe-artifacts are quite prominent in the final result and uncertainty estimate here. In addition, it would be nice if the uncertainty estimate made some attempt to account for the uncertainty in the input datasets, which have uneven spatial resolution in Antarctica even for variables that are relatively well-constrained like seismic velocity or Curie Depth.

405 We thank the reviewer for these comments. As described in our response to the Major Comments, we have addressed this issue by projecting the data from EPSG:4326 to EPSG:3031 (Antarctic Polar Stereographic Projection) before interpolation. Regarding uncertainty estimation, we acknowledge that the uncertainty in input datasets, particularly their uneven spatial resolution in Antarctica even for relatively well-constrained variables like seismic velocity or Curie depth, contributes to the overall aleatoric uncertainty captured in our framework. And we have substantially improved our uncertainty quantification by implementing a Bayesian framework that decomposes total predictive uncertainty into aleatoric and epistemic components:

410 *“The distribution of uncertainty components (Fig. 8a) reveals that aleatoric uncertainty constitutes the dominant fraction of total uncertainty throughout the study region. The median aleatoric uncertainty ($\sim 650 \text{ mW}^2 \text{ m}^{-4}$) substantially exceeds the median epistemic uncertainty ($\sim 70 \text{ mW}^2 \text{ m}^{-4}$), indicating that inherent variability in heat flow observations and unresolved local geological heterogeneity—including the uneven spatial resolution of input datasets in Antarctica—represent the primary sources of predictive uncertainty. This finding suggests that while our model has sufficient capacity*

415 *and training data to capture the underlying patterns, the irreducible observational noise and small-scale geological complexity impose fundamental limits on prediction accuracy. ”*

Figure 7

420 It would be better to show signed difference rather than absolute difference here. It is important to know which estimate is hotter! This would be a good place to use the blue-white-red color scale from figure 1.

In addition, there are quite a few additional published estimates that you could compare your model against. Additional comparison datasets include: Shapiro and Ritzwoller (2004); Maule et al., (2005); Purucker et al., (2012); An et al., (2015); Lucazeau, (2019), Haeger et al., (2022); Hazzard and Richards, (2024).

425 We thank the reviewer for these suggestions. We have revised Figure 7 to show signed differences rather than absolute differences, using a blue-white-red diverging color scale to clearly indicate which estimate predicts higher or lower GHF values relative to our model. We have expanded our comparison to include additional published estimates:

430 *“To quantitatively assess the relationship between our predictions and existing models, we computed the spatial differences between our GHF map and six published estimates: Fox Maule et al. (2005), Martos et al. (2017), Shen et al. (2020), Lösing and Ebbing (2021), Stål et al. (2022), and Hazzard and Richards (2024). ”*

L356: “In instance...”

435 Should be, “For instance...”

Thanks. Did it!

Figure 8

440 It appears that many of the observations that you use to validate your model are actually located on the seafloor around Antarctica. While it certainly makes sense to include these data points when so few in situ observations are available, does it really make sense to compare your model against these data when you excluded marine observations from your training data?

445 Indeed, it is inconsistent to validate our model against marine observations when we excluded marine data from our training dataset. We have removed this validation section from the revised manuscript.

L389-392: “This discrepancy may result from the heterogeneity of local geologic features, differences in raw data processing methods, or the influence of complex processes such as shallow water circulation and unsteady convection in the lithosphere, and further studies are needed to elucidate the underlying mechanisms.”

In addition, the discrepancy between your results and those of Shroeder et al. (2014) could be the result of model assumptions made by Schroeder et al. They made very specific and potentially limiting assumptions about the form of the subglacial hydrological system when constructing their inverse model, and those assumptions could potentially introduce errors into their result.

L389-392: We thank the reviewer for this insightful comment regarding the potential influence of model assumptions in Schroeder et al. (2014). We have substantially revised the Discussion section, and this paragraph has been removed as part of the reorganization. Please refer to the revised Discussion for details.

Section 6 Data Availability

This section should be after the Conclusions section, not before it.

Thanks. We have moved the Data Availability section to follow the Conclusions section.

L426: Zenodo link

It would be nice if this link also contained the processed and gridded datasets used as input to your model. While it is true that these datasets are all available at their original sources, it would be nice if it were possible for interested users to access the gridded inputs that you created for your model at one place.

We agree. We have updated the Zenodo repository to include all processed and gridded input datasets used in our model, in addition to the final GHF predictions and code.

L437: “...which is consistent with the active geological structures.” Rephrase this, this sounds awkward. Perhaps try: “...which is consistent with the locations of present-day tectonic and volcanic activity.”

We agree that the suggested phrasing is clearer and more precise. We have substantially revised the Discussion section, and this paragraph has been removed as part of the reorganization. Please refer to the revised Discussion for details.

References: The references should be in alphabetical order, not in citation order.

Thank you for pointing this out. We have reorganized the reference list into alphabetical order by first author's surname, following standard formatting conventions.