Replies to Referees' comments on "Drawing lessons for multi-model ensemble design from emulator experiments: application to future sea level contribution of the Greenland ice sheet" (egusphere-2025-52)

We would like to thank both Referees for the constructive comments. We agree with most of the suggestions and, therefore, we have modified the manuscript to take on board the comments and suggestions. We recall the reviews and we reply to each of the comments in turn (outlined in blue). The page and line numbers are those of the document with track-changes.

Additional changes

Since the submission of this manuscript, we had the opportunity to include HIRHAM RCM in the MME as well. In addition to the modifications suggested, we have thus included this third RCM model in the analysis. A major advantage is to define an experiment where only the members of HIRHAM and RACMO are used ('woMAR' experiment) in addition to the 'MAR only' experiment. We expect that this modification brings new insights and strengthens our conclusions which are, in the original version of the manuscript, based on two RCMs only. The results discussed in Sect. 4 have been modified accordingly.

Referee #1:

Review of "Drawing lessons for multi-model ensemble design from emulator experiments: application to future sea level contribution of the Greenland ice sheet" by Rohmer et al. The paper "Drawing lessons for multi-model ensemble design from emulator experiments: application to future sea level contribution of the Greenland ice sheet" investigates the impact of different choices in constructing an emulator capable of predicting the future contribution of the Greenland Ice Sheet to sea-level rise. Specifically, the authors examine how the selection of models and parameters associated with different physical processes and scales (climate scenarios, regional climate models, ice-flow models, and ice-flow parameters) influences the emulator in terms of its fidelity, the estimated contribution to sea-level rise, and the associated uncertainty.

This study represents a valuable contribution to our understanding of multi-model ensemble approaches and emulator design. The numerical experiments are clearly described, and the results are noteworthy. Overall, the paper provides an interesting addition to the scientific literature. Nonetheless, I have a series of comments and questions that I would like the authors to address. On the one hand, in terms of the way the manuscript is written, some sections could be streamlined for conciseness (see my first general comment below, as well as specific comments). On the other hand, I have a series of more fundamental general comments on key aspects of the manuscript. There are three main comments, described herefater, which are followed by a series of specific comments addressing more detailed points. Once these concerns have been satisfactorily addressed, I will be happy to recommend the paper for publication in the Special Issue: Improving the contribution of the land cryosphere to sea level rise projections.

We thank Referee #1 for the positive analysis. We have taken into account the comments and suggestions. In what follows, we describe in detail the corrections.

General comment 1

This paper lies at the intersection of two fields: glaciological modeling and statistical methods. Such interdisciplinary studies are particularly valuable, as the glaciology community may not be fully familiar with statistical techniques, while statisticians may not always be aware of the challenges involved in estimating future sea-level rise. Furthermore, comparative studies have gained increasing importance in glaciology, as they help assess the robustness of different modeling approaches. Given this, it is crucial to also investigate how these comparisons are constructed in the first place. In this context, the present study is highly relevant.

That being said, I believe the paper could be more explicit about its practical conclusions, specifically conclusions (1) and (2) mentioned in the abstract. I would also suggest highlighting the conclusions related to the (different) impact of the RCM/ISM choice in the abstract and in the introduction (see lines 313–315).

Additionally, some technical details could be either removed or moved to an appendix or supplementary materials. The reasoning behind these suggestions is that most readers of this journal are likely geoscientists, who will primarily be interested in the study's results. By streamlining technical details, the paper's key findings —which are noteworthy— could be emphasized more effectively, improving its accessibility.

We thank Referee #1 for these suggestions. We totally agree that the results described in our manuscript should be transferred to a wide readership that is not necessarily specialists of the methods.

To do so, the following modifications have been made:

- The technical details of the emulator implementation are placed in Appendix A. We also added a new Appendix D to detail the methods for the performance analysis suggested by Referee #2. Finally, the specific comments of Referee #2 are also integrated in Appendices to decrease the level of technicality of the core text;
- The conclusions of the RCM/ISM are more highlighted in the abstract and in the conclusions:
- In order to improve the transferability of our message, some choices made for the representation of the results are modified to be more consistent with IPCC standards, i.e. with more largely shared practices:
 - We slightly change the definition of GSAT by computing the difference between the temperature at the considered year and the mean temperature over the period 1995-2014;
 - We analyse in Sect. 3.3 the changes in the likely range, i.e. 17th and 83rd percentile instead of the 5th and 95th percentile. We believe that the use of the IPCC calibrated language can also ease the communication of our results.

I also wonder whether subsections 2.4 and 3.2 could be simplified by focusing on just one of the two indicators (DS and Dh), e.g. DS. While some variability has been observed in the results, these indicators are strongly correlated (see figure 6). Furthermore, from a practical point-of-view, what really matters is the number of runs to be made (i.e., of DS). Similarly, figure captions could be streamlined by omitting details that may not be particularly useful (e.g., the exact method used for kernel density estimation; see specific comments below). The technical lines 196–219 could also be relegated to supplementary materials.

We agree with Referee #1's comment: the number of members is a key criterion, and the evolution of the emulator's predictive capability clearly highlights this aspect. Our primary idea for proposing a second criterion was to be able to reflect how much the distribution of the members is modified in addition to the size. As rightly indicated by Referee #1, our second

criterion D_h fails, however, to highlight key situations because of the too strong correlation with D_S . This goes also in the same line with Referee #2's comment (minor comment 2.1).

Therefore, we propose to remove this analysis from the main text. We propose then to support the discussion in Sect 4.1 with a complementary analysis (Supplementary materials S3) based on a well-established, and more widely used, criterion for comparing different probability distributions, namely the Kolomogorov-Smirnov (KS) criterion, instead of a criterion constructed from 'scratch'.

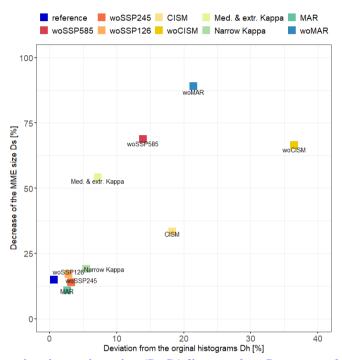


Figure S7: Position of the emulator's experiment in a (D_h, D_s) diagram where D_s measures the relative decrease in the MME size after applying the experiment, and D_h measures the deviation of the histograms from the original ones (see Supplementary Material S3). The blue-coloured marker refers to the reference solution defined as the mean value over the 25 iterations of the random validation exercise, described in Sect. 2.4, applied to the original dataset.

General comment 2

Regarding the paper's methodology, I have some questions about how the way ensembles are introduced. Specifically, lines 34–36 state: "Each member of a MME should evenly span a representative and exhaustive set of plausible realizations of the combined sources of uncertainty, (...), equally represented by a single model run". This suggests that each member of the ensemble should have the same weight.

However, I find this somewhat misleading, and I believe the authors could elaborate further on this choice. The assumption that all runs should have equal weight holds only if our current knowledge suggests they are equally probable. From a Bayesian perspective, this would correspond to assuming a uniform prior distribution. However, this assumption may not always be justified. For example:

- (i) Runs from lower-resolution models might be considered less reliable than those from higherresolution models as they might not capture relevant small-scale processes.
- (ii) Some values of the uncertain parameters might be less probable if they lead to results that deviate significantly from current observations.

Formally, these concerns can be addressed by updating the weights of each run based on their likelihood given observational data (e.g., Aschwanden and Brinkerhoff, 2022; Nias et al., 2023).

I understand that the authors did not include such a step in their analysis, as their focus was on assessing the emulator's capabilities. However, it seems to me that this point should be discussed in the methodology or discussion section for two reasons. First, to clarify for the reader that the choice of equal probability stems from an assumption about our current knowledge and that alternative approaches are possible. Second, because weighting model runs based on observational constraints is an emerging direction in the field, and this should be discussed in the context of future ISMIP ensemble designs. This could also be mentioned as a perspective for future work, as it would be interesting to see whether the conclusions remain valid when runs are weighted as a result of a calibration.

We agree with Referee #1 that this point merits further discussion, particularly in view of the forthcoming ISMIP7. The primary goal our study was to study the influence of different factors for a given MME, i.e. discovering the influence of groups of members. Thus, the implicit assumption of this procedure is that we do not assign any weight to the members *a priori*. This is clarified on page 8, in Sect. 2.3 (line 180); this was also Referee #2's suggestion (Minor comment 4(b)).

This clarification is more particularly useful in Sect. 4 regarding the meaning of the emulator-based probabilistic predictions. These projections do neither represent calibrated uncertainty accounting for model-observation misfits nor the *slc* probability distribution from the MME, because the uniform distribution over the input space is not representative of the MME itself (e.g., the minimum spatial resolution is clearly not uniform between 1 and 40 km, see Fig. 3). This is clarified on page 11, in Sect. 2.4.2 (lines 227-234).

In addition, the alternative option based on weighting either through expertise (as illustrated by Referee #1 with the resolution) or based on model-observation misfits is now discussed in Sect; 5. We propose to highlight the benefits of the weighting approach, but also the challenges to do it, namely: (1) the need for good quality data; (2) the need for data over a large period in the past; (3) the need for some types of ISMs to adjust or adapt their implementation. This is clarified on page 20, in Sect. 5 (lines 424-426).

General comment 3

My third general remark concerns the parameter κ associated with the calving rate. It is unclear why this particular parameter was chosen over others. From reading the paper, the rationale behind this choice is not obvious—perhaps it is based on modeling considerations or supported by previous studies? If so, it seems to me that the authors should provide a stronger justification for including this specific parameter.

More fundamentally, I wonder whether comparing the effect of κ to that of RCP scenarios or RCMs is entirely meaningful. This comparison contrasts the impact of a single parameter of an ice-flow model with that of an entire climate scenario or a regional climate model, which incorporate numerous physical parameters. Given this, it is perhaps unsurprising that the effect of κ appears quite limited.

We thank Referee #1 for this suggestion. We feel that some clarification about κ should here be given.

In this study, we rely on a standard approach for integrating ocean forcing, i.e. based on an empirically derived retreat parameterization for tidewater glaciers (Slater et al., 2019, 2020) that is forced by a RCM-based run-off and ocean temperature changes in seven drainage basins around Greenland. In this implementation, retreat and advance of marine-terminating outlet glaciers in the ISMs are prescribed as a yearly series of maximum ice front positions (Nowicki et al., 2020).

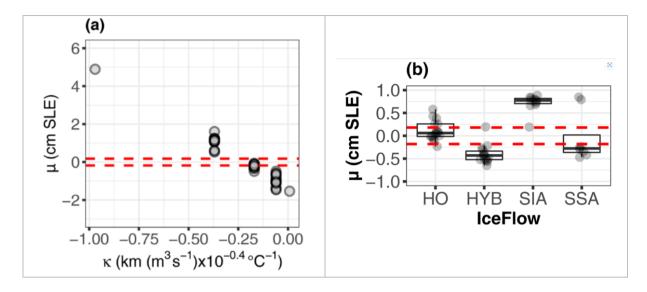
Here, κ is not thought as a parameter of the ice-flow model, it rather represents the sensitivity of the ocean forcing as a whole. It may be thought of as defining the sensitivity of the downscaling from global model to local ice sheet scale, similar to the combined parameter choices in RCMs for downscaling climate conditions. In the studied MME, we have different RCMs, which have different sensitivities and produce different melt for the same global forcing. Since we have only one approach to 'downscale' the ocean forcing, κ is sampling that uncertainty in a similar way.

We recognize that the κ -based approach remains a strong simplification of the complex interaction between marine-terminating outlet glaciers and the ocean, for which physically based solutions are in development but not available for all models. However, it should be underlined that the advantage of this retreat parameterization is to be applicable in the wide variety of models under consideration. Furthermore, it is currently the most widely used approach for producing large ensemble for sea level projections, as done for instance by Edwards et al. (2021).

This is clarified on pages 3-4, in Sect. 2.1 (lines 96-100).

To illustrate this, one could consider a similar comparison in the opposite direction: assessing the impact of choosing a glaciological model of varying complexity (e.g., full Stokes, BP, or SIA) against a single parameter from a RCM. This would likely lead to the conclusion that the specific parameter from the RCM has a minimal influence. Therefore, I wonder whether including κ as an isolated parameter in this study is fully justified. Could the authors maybe clarify its relevance within the broader context of the study's objectives?

In response to the specific comment of Referee #1, our previous study (Rohmer et al., 2022) highlighted the high importance of κ compared to other uncertainties, in particular some of them related to the complexity of the numerical model as suggested by Referee #1, i.e. the choice in the numerical method (Finite Difference, Finite Element), the grid resolution, and the ice flow formulation (approximation, higher order, hybrid). To illustrate, the following figure (adapted from Fig. 7 and Fig. 8 of Rohmer et al., (2022) based on ISMIP6 MME) shows the sensitivity index (denoted μ) that measures the contribution, in terms of sea level equivalent SLE, depending on the value of κ (panel a) or of the choice in the ice flow formulation (panel b). The influence measured by μ for κ is on the order of 1-2 cm (at most 5cm) whereas it remains on the order of 0.5-1cm for the ice flow method, hence confirming a large importance of this parameter.



References

Edwards, T. L., Nowicki, S., Marzeion, B., Hock, R., Goelzer, H., Seroussi, H., et al. (2021). Projected land ice contributions to twenty-first-century sea level rise. Nature, 593(7857), 74-82.

Nowicki, S., Goelzer, H., Seroussi, H., Payne, A. J., Lipscomb, W. H., Abe-Ouchi, A., Agosta, C., Alexander, P., Asay-Davis, X. S., Barthel, A., Bracegirdle, T. J., Cullather, R., Felikson, D., Fettweis, X., Gregory, J. M., Hattermann, T., Jourdain, N. C., Kuipers Munneke, P., Larour, E., Little, C. M., Morlighem, M., Nias, I., Shepherd, A., Simon, E., Slater, D., Smith, R. S., Straneo, F., Trusel, L. D., van den Broeke, M. R., and van de Wal, R.: Experimental protocol for sea level projections from ISMIP6 stand-alone ice sheet models, The Cryosphere, 14, 2331–2368, https://doi.org/10.5194/tc-14-2331-2020, 2020

Rohmer, J., Thieblemont, R., Le Cozannet, G., Goelzer, H., and Durand, G.: Improving interpretation of sea-level projections through a machine-learning-based local explanation approach, Cryosphere, 16, 4637–4657, https://doi.org/10.5194/tc-16-4637-2022, 2022.

Slater, D. A., Felikson, D., Straneo, F., Goelzer, H., Little, C. M., Morlighem, M., Fettweis, X., and Nowicki, S.: Twentyfirst century ocean forcing of the Greenland ice sheet for modelling of sea level contribution, The Cryosphere, 14, 985–1008, https://doi.org/10.5194/tc-14-985-2020, 2020.

Slater, D. A., Straneo, F., Felikson, D., Little, C. M., Goelzer, H., Fettweis, X., and Holte, J.: Estimating Greenland tidewater glacier retreat driven by submarine melting, The Cryosphere, 13, 2489–2509, https://doi.org/10.5194/tc-13-2489-2019, 2019.

Specific Comments

We thank Referee #1 for the specific comments. The revised version of the manuscript now incorporates all of them.

(1) [Line 12] 'projection and the quantification of its uncertainty' \rightarrow 'projections and the quantification of their uncertainties'.

This has been corrected.

(2) [Lines 15, 17, and 65] You use 'experiments' for two distinct concepts: numerical simulations (e.g., line 16) and numerical tests (e.g., line 15). Consider using separate words to avoid any confusion.

To avoid confusion, we now refer to the second type of 'experiments' as 'members'.

(3) [Line 16] '(Regional Climate Model RCM, or Ice Sheet Model ISM)' → '(Regional Climate Model; RCM, or Ice Sheet Model; ISM)'.

This has been corrected.

- (4) [Line 19] Consider removing 'utmost' as it might be overly formal. This has been replaced by 'high'.
- (5) [Line 25] 'projection and the quantification of its uncertainty' \rightarrow 'projections and the quantification of their uncertainties'. This has been corrected.
- (6) [Line 26] 'co-ordinated sets of numerical experiments'→'sets of numerical experiments'. This has been corrected.
- (7) [Line 47] Consider adding references related to (machine-learning-based) emulators. A list of real case applications is now provided on page 2 in Sect. 1 (lines 51-52) as follows: "[...] like linear-regression (Levermann et al., 2020), Gaussian process regression (Edwards et al., 2021), random forest regression (Rohmer et al., 2022), and deep learning-based methods (Van Katwyk et al., 2025)".

Added references

Levermann, A., Winkelmann, R., Albrecht, T., Goelzer, H., Golledge, N. R., Greve, R., Huybrechts, P., Jordan, J., Leguy, G., Martin, D., et al.: Projecting Antarctica's contribution to future sea level rise from basal ice shelf melt using linear response functions of 16 ice sheet 600 models (LARMIP-2), Earth System Dynamics, 11, 35–76, https://doi.org/10.1175/JCLI-D-23-0580.1, 2020.

Van Katwyk, P., Fox-Kemper, B., Nowicki, S., Seroussi, H., & Bergen, K. J. (2025). ISEFlow v1. 0: A Flow-Based Neural Network Emulator for Improved Sea Level Projections and Uncertainty Quantification. EGUsphere, 2025, 1-32.

- (8) [Line 47] I might be a bit picky here, but I would argue that the key advantage of statistical emulators is their low computational cost; being able to predict the model response at untried input values is only useful if it can be done at a reasonable cost.

 We totally agree with Referee #1. We have now underlined this aspect.
- (9) [Line 63] Please be consistent with your use of acronyms: either your define what you mean by RCM and GCM, or you use directly the corresponding acronyms. Also, a table of acronyms would be useful in the paper.

A table has been added in Appendix E.

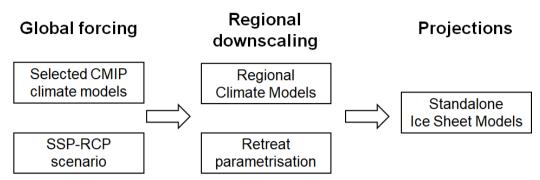
- (10) [Line 63] Avoid using 'validation tests' as this can lead to confusion when it comes to glaciological modeling (for which 'validation' has another meaning).

 We now use the term "numerical experiments".
- (11) [Line 76] '(Goelzer et al. (2020): in particular (...))' \rightarrow '(Goelzer et al., 2020; in particular (...))'.

This has been corrected.

(12) [Line 79] Consider adding a schematic displaying the modeling chain and indicating where modeling choices (MME inputs) are introduced. This could be very useful to effectively obtain an overview of the context.

A new Figure 1 has been added to clarify the workflow.



New Figure 1: General forcing approach for Greenland ice sheet model projections. The questions relevant for the MME design (detailed in Table 2) are related to the modelling choices made for each of the boxes.

(13) [Line 80] Here you define again what a RCM is. If you have already defined it before that is not necessary.

We have removed this part.

- (14) [Table 1] Ensure consistent formatting of symbols (italics vs. non-italics). This has been corrected.
- (15) [Table 1] Consider renaming 'Symbol' to 'Symbol/Acronym' or simply 'Name' to clarify that most entries are acronyms.

The term 'Name' is now used.

- (16) [Table 1] 'Sliding basal law'→'Sliding law' or 'Basal friction law'. The second term is now used.
- (17) [Line 101] Consider defining 'input setting' explicitly, e.g. as a particular combination of inputs.

This has been specified as suggested.

Referee #2:

This study develops a Random Forest (RF) emulator to emulate Greenland 2100 sea level contribution (slc) output from a Multi-model ensemble (MME). In particular, the RF is trained using a set of 7 inputs, associated with the climate scenario, the ice sheet model (ISM) used, the regional climate model (RCM) used, and different settings of the ISM run. The authors investigate how changing the MME design leads to changes in the emulator performance and in its range of emulated slc. Based on these metrics, they provide guidelines for future MME designs that aim at estimating future ranges of slc from the Greenland ice sheet.

This study addresses an important and difficult question: how can we improve the design of MMEs to provide the best information about the probability density function (PDF) of future slc? The concept underlying this study is that the MME itself does not need to characterize this PDF, but that it should be designed optimally such that an emulator can do this characterization a posteriori. This is a valid and efficient approach to uncertainty quantification. It is also a challenging topic, and work on this topic is important. However, at this stage, I believe that several points need to be improved to make this study a valuable contribution in addressing this question. The authors make recommendations and they "expect these recommendations to be informative for the design of next generations of MME" (L22). But I believe that their recommendations are dependent on many assumptions or choices that they made, without always justifying them and making them clear to the reader. Furthermore, more methodological details about the RF emulator are needed because all the results presented depend on the emulation and, therefore, the RF design influences strongly any interpretation, and thorough RF evaluation is critical as well. I detail my concerns in this review, which consists of Major and Minor comments. I do not provide technical comments at this stage of the reviewing process because I believe that the more substantial aspects should be addressed first. Line numbers in this review correspond to the preprint manuscript.

We thank Referee #2 for the in-depth analysis of the manuscript. In what follows, we provide details of the corrections made.

Major comment 1: Inherent assumptions associated with the MME

Many of the conclusions are strongly dependent on the particular MME used in this study. I have several reservations about this.

First, it is unclear to me how the MME used in this study was acquired and designed. The only details provided about the MME are (L74): "We focus on the sea level contribution from the Greenland ice sheet (GrIS) in 2100 based on a new MME study performed for the European Union's Horizon 2020 project PROTECT (http://protectslr.eu). Some modelling choices are taken from the protocols of the ISMIP6 initiative (Goelzer et al. (2020): in particular, the two main emissions scenarios, and the main model parameter explored."

Has this MME been peer-reviewed? Why are the authors not using the well-established ISMIP6 MME? The latter MME also has the advantage of providing a larger set of experiments, notably including many more ISMs than the MME of this study. At least, why has the MME not been combined with the ISMIP6 MME? Also, given that no publication describing the MME is referenced, I believe that it is important to give many more details about the MME configuration: Did all ISMs run under high- and low-warming forcing? Are the 15 global climate models used in the MME well-balanced across the runs? Etc.

The MME has been peer reviewed within the H2020 Protect project and has been submitted to Special Issue of the Cryosphere journal. In the revised version of the manuscript, we refer to the egusphere preprint for further details (Goelzer et al., 2025). Since the release on the egusphere platform takes some time, Referee #2 is invited to refer to:

https://drive.google.com/file/d/1S00IHWGa34mLNlLOyrHEN6Sp7q6-

EbLY/view?usp=sharing

** Please to do not distribute / use outside the scope of this review**

Reference added

Goelzer, H., Berends, C. J., Boberg, F., van den Broeke, M., Durand, G., Edwards, T., Fettweis, X., Gillet-Chaulet, F., Glaude, Q., Huybrechts, P., Le clec'h, S., Mottram, R., Noel, B., Olesen, M., Rahlves, C., Rohmer, J., van de Wal, R. S. W.: Extending the range and reach of physically-based Greenland ice sheet sea-level projections. Preprint egusphere-2025-3098, 2025.

To specifically reply to Referee #2 comment, we confirm that the experimental design builds on the ISMIP6 protocol with four different ice sheet models and extends it to more fully account for uncertainties in sea-level projections for the GrIS. This is underlined in Sect. 2. The Protect MME is thought as an extension of ISMIP6 MME regarding different aspects:

- we have included a wider range of CMIP6 climate model output, more climate change scenarios (SSP126, SSP245, SSP585);
- we have provided retreat forcing before 2015 that is calculated from reconstructions of past runoff and ocean thermal forcing. This allows for a consistent forcing of the models in past and future and to consider historical retreat of the outlet glaciers, which was an important source of mass loss after 1990;
- we have provided surface mass balance forcing from several RCMs, i.e. MAR and RACMO for the MME used for this study.

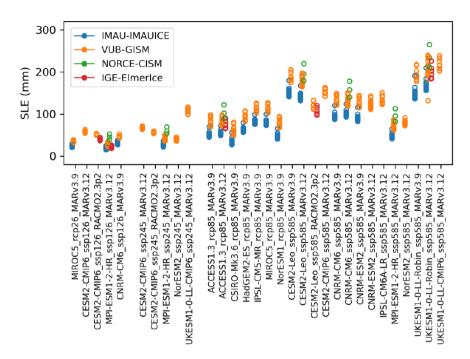
Since the submission of this manuscript, we had the opportunity to include HIRHAM as well. We believe that adding this third RCM model brings new insights and makes our conclusions more robust (based on two RCMs in the original version of the study). The results discussed in Sect. 4 have been modified.

To appreciate this extension, the following table gives the repartition of the members with respect to the RCP/SSP scenarios for ISMIP6 MME and for Protect MME (before inclusion of HIRHAM). In addition, it should be noted that the total number of members has increased by a factor of 4 compared to ISMIP6 MME.

	ISMIP6 MME	Protect MME
RCP26	23	40
RCP85	156	319
SSP126	18	189
SSP245	0	189
SSP585	59	566
Total number	256	1303

As a second illustration, the following figure gives an overview of the results of the ensemble of projections at the year 2100 for all available Earth System models (ESMs), RCMs and ISMs under high, med and low retreat sensitivity. This allows to appreciate, graphically, how well the range of sea level is covered. It should however be underlined that, although the collection

of forcing data covers a wide range of variations across different ESMs and scenarios, it ultimately still represents an 'ensemble of opportunity' similarly as for ISMIP6 MME. Our study aims to address the potential implications of this characteristic.

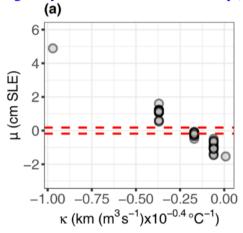


Second, the conclusions of this study are strongly dependent on the initial MME used in the emulation process. For example, the authors argue that there is "a quasi-linear relationship between κ and slc" (L306). But this conclusion is based only on the set of 4 ISMs used in this study: CISM, IMAUICE, GISM, and ElmerIce. Furthermore, given that very little experiments were performed by ElmerIce and GISM, I assume (I need to assume here because no information is provided on the design of the MME) that these two models may well have only been run with a single κ value. In this case, the "quasi-linear relationship" would be derived only from two ISMs. Given that different ISMs can show very different sensitivities to movement of the tidewater glacier front and grounding line positions, this conclusion could well be very different if other ISMs are included. So, if one was to perform a similar study with the ISMIP6 MME, would the "recommendations" for future MME design be different? As another example, I mention above that only 4 ISMs are included in the MME, two of which account for > 90% of the simulations. By excluding CISM from the training experiments, the authors then make "recommendations" about the ability of the emulator to estimate the slc simulated by ISMs not included in the design. Here also, this evaluation depends critically on how similar simulated slc from CISM is to the simulated slc from IMAUICE. This similarity depends on numerous aspects that are specific to these two particular models. I would expect that the "recommendations" would be very different if other ISMs (ISSM, PISM ...) show more or less similarity with CISM.

We thank Referee #2 for the insightful analysis. Here we feel that some clarifications about κ should be given. As described in Sect. 2, on page 3-4 (lines 96-99), we rely on a standard approach for integrating ocean forcing, i.e. based on an empirically derived retreat parameterization for tidewater glaciers (Slater et al., 2019, 2020). In this approach, κ is <u>not</u> a parameter in the ice-flow model; it rather represents the sensitivity of the ocean forcing as a whole. It may be thought of as defining the sensitivity of the downscaling from global model to local ice sheet scale, similar to the combined parameter choices in RCMs for downscaling

climate conditions. In the studied MME, we have different RCMs, which have different sensitivities and produce different melt for the same global forcing.

Regarding the specific comment on the quasi linear behaviour, we expect this relationship not to change much by adding more ISMs because κ is "external" to the ISM as afore-described. To support this result, Referee #2 should refer to the following figure adapted from Fig. 7 of Rohmer et al. (2022) based on ISMIP6 MME: it shows the sensitivity index (denoted μ) that measures the contribution, in terms of sea level equivalent SLE, depending on the value of κ . A quasi-linear trend has here been identified. To complement this analysis, we rely on Supplementary Materials S3, which provides details on the quasi-linear relationship of the partial dependence plot (as originally described in the manuscript).



Finally, we agree with Referee #2 that a more careful attention should be paid not to 'over-interpret' our results by making the recommendations too general. The conclusions are nuanced and reformulated in this sense. In addition, we also propose to modify the title to highlight that our results are linked to the specificities of our ensemble as follows: "Lessons for multi-model ensemble design from emulator experiments: application to a large ensemble for future sea level contributions of the Greenland ice sheet".

Reference

Rohmer, J., Thieblemont, R., Le Cozannet, G., Goelzer, H., and Durand, G.: Improving interpretation of sea-level projections through a machine-learning-based local explanation approach, Cryosphere, 16, 4637–4657, https://doi.org/10.5194/tc-16-4637-2022, 2022.

Such assumptions are not made explicit by the authors. This could very well be misleading to the readership targeted by the authors, especially those less familiar with ice sheet modeling (e.g., "stakeholders" (L328) and "coastal adaptation practitioners" (L332)).

We totally agree with Referee #2. The clarification on κ has been added in Sect. 2 on page 3-4 (lines 96-99).

Major comment 2: Characterization of uncertainty

The authors use their random forest (RF) emulator such that "changes in the emulator's predictive performance and the emulator-based probabilistic projections provided information on several aspects" (L18). After reviewing the manuscript, I identify remaining limitations about the RF emulator regarding uncertainty characterization.

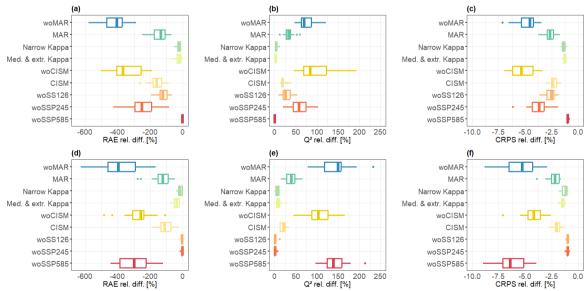
The authors use changes in the predictive performance of the RF as a proxy for uncertainty remaining about a hypothetical MME (here, a MME excluding some of the experiments). But this metric is sensitive to the particular machine learning model used for the emulation. Here, the emulation output is thus conditioned on the RF architecture, with a single fixed combination of hyperparameters. Is any decrease in predictive performance of this specific RF therefore a meaningful assessment of uncertainty imputable to the MME design? This question is critical, because the conclusions of this study use this as a fundamental assumption.

This issue is further exacerbated by the fact that the RF does not provide probabilistic output. By this, I mean that the RF only provides a point estimate. There is no uncertainty quantification. Ideally, the design of a MME should target the strongest reduction in posterior covariance (i.e., the uncertainty remaining given the current MME). But this particular RF emulator does not provide such metric. This could be addressed by choosing another architecture (e.g., Gaussian processes, Williams and Rasmussen (2006)), by subsampling techniques for RF models (Mentch and Hooker, 2016), or by adapting the RF to output conditional quantiles (Meinshausen and Ridgeway, 2006).

We thank Referee #2 for this suggestion. We agree that complementing the study with uncertainty quantification of the emulators itself bring new insights and allow us to better discuss the results. As suggested, we propose to implement the quantile random forest emulator, denoted qRF, for both the experiments on the emulator's performance (Sect. 3.2) and for the probabilistic projections (Sect. 3.3).

For the former application, the quantile random forest provides estimates of quantiles at any order τ , denoted $q^{\tau}(slc|\mathbf{x}^*)$ for a given instance of the input variables \mathbf{x}^* . The quantiles can directly be used to define the prediction intervals at any level α : $[q^{(1-\alpha)/2}(slc|\mathbf{x}^*);q^{(\alpha+1)/2}(slc|\mathbf{x}^*)]$.

The new Figure 7 allows to verify the levels of predictability for different intervals of GSAT changes.



New Figure 7: Relative difference (in %) of the performance criteria considering the lowest GSAT values below 2.14°C (top) and the highest GSAT values above 3.83°C (bottom) for RAE (a, d), Q^2 (b, e), and CRPS (c, f).

When performing the probabilistic predictions (Sect. 2.4.2), the emulator uncertainty is propagated in addition to the uncertainty of the different input variables based on the following procedure:

(Step 1) Draw N random realisations of the input variables $\tilde{\mathbf{x}}$;

(Step 2.1) Draw N random number \tilde{u} between 0 and 1 by assuming a uniform random distribution;

(Step 2.2) Compute the N values $\widetilde{slc} = q^{\widetilde{u}}(slc|\widetilde{\mathbf{x}})$ given \widetilde{u} and $\widetilde{\mathbf{x}}$ using the qRF model;

(Step 2.3) Compute the quantile $Q_{\widetilde{u}}^{\alpha}$ at the chosen level α from the set of N values of \widetilde{slc} ;

(Step 3) Repeat n times Steps 2.1 to 2.3. At Step 2.2, \widetilde{slc} are calculated for the same set of random input variables $\tilde{\mathbf{x}}$ defined at Step 1, but for a newly randomly generated set of levels \tilde{u} defined at Step 2.1. At Step 2.3, the newly calculated quantiles $Q_{\tilde{u}}^{\alpha}$ vary at each of the repetitions, since each time, new random levels \tilde{u} of the qRF conditional quantiles are generated at Step 2.1.

The output of the procedure is a set of n quantile values $(Q_{\widetilde{u}^{(1)}}^{\alpha}, Q_{\widetilde{u}^{(2)}}^{\alpha}, \dots, Q_{\widetilde{u}^{(n)}}^{\alpha})$. The variability of the set reflects the emulator uncertainty and can be summarized by the τ % confidence interval with lower and upper bounds defined by the $(1-\tau)/2$, and the $(1+\tau)/2$ quantile of $Q_{\widetilde{u}}^{\alpha}$. In this study, we choose N=10,000, n=100 and $\tau=90\%$.

In addition, we propose to add a new performance indicator to analyse the changes in the emulator's performance in terms of reliability of the predictive probabilistic distribution. This is done using the continuous ranked probability score, denoted *CRPS*, as used for validating probabilistic weather forecast (Gneiting et al., 2005). To evaluate the *CRPS* score, the formulation based on quantiles (Berrisch and Ziel (2024): Eq. 2) is used:

$$CRPS = 2 \int_0^1 B(q^\tau(slc|\mathbf{x}^*), slc^{true}) \, \mathrm{d}\, \tau \approx \frac{2}{P} \sum_{\tau \in \Gamma} B(q^\tau(slc|\mathbf{x}^*), slc^{true})$$
 where the term $B(q^\tau(slc|\mathbf{x}^*), slc^{true})$ is defined as
$$\{(1-\tau)(q^\tau(slc|\mathbf{x}^*) - slc^{true}) \text{ if } slc^{true} < q^\tau(slc|\mathbf{x}^*) \\ \tau(slc^{true} - q^\tau(slc|\mathbf{x}^*)) \text{ if } slc^{true} \geq q^\tau(slc|\mathbf{x}^*) \\ \tau(slc^{true} - q^\tau(slc|\mathbf{x}^*)) \text{ if } slc^{true} \geq q^\tau(slc|\mathbf{x}^*) \\ \text{sea level contribution, and where the quantiles } q^\tau(slc|\mathbf{x}^*) \text{ are evaluated using the trained qRF} \\ \text{model at given instance of the input variables } \mathbf{x}^* \text{ for an equidistant dense grid of quantile levels} \\ (\tau_1, \dots, \tau_P) \text{ with } \tau_i < \tau_{i+1} \text{ and } \tau_{i+1} - \tau_i = 1/P. \text{ In this study, we consider level } \tau_1 = 5\% \text{ and } \\ \tau_P = 95\% \text{ with } 1/P = 5\%.$$

This score jointly quantifies the calibration of qRF probability distribution, i.e. the reliability of the estimation, and its sharpness (i.e. the concentration/dispersion of the probability distribution). The lower *CRPS*, the higher the quality of the qRF probabilistic predictions, with a lower limit of zero.

Finally, we underline in Sect. 5 the problem of model uncertainty related to the construction of the emulator; in particular the problem of hyperparameters' tuning and its relatively lesser impact for random forest models (Bischl et al., 2023).

Added references

Berrisch, J., Ziel, F., 2024. Multivariate probabilistic crps learning with an application to dayahead electricity prices. International Journal of Forecasting, 40(4), 1568-1586, doi:10.1016/j.ijforecast.2024.01.005

Bischl, B., Binder, M., Lang, M., Pielok, T., Richter, J., Coors, S., et al., 2023. Hyperparameter optimization: Foundations, algorithms, best practices, and open challenges. Wiley

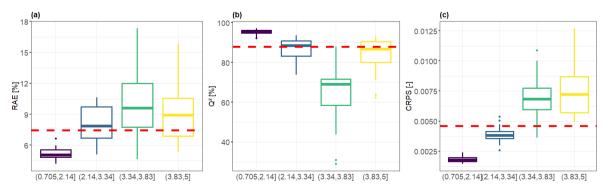
Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 13(2), e1484, doi: 10.1002/widm.1484

Gneiting T, Raftery AE, Westveld III AH, Goldman T. 2005. Calibrated probabilistic forecasting using ensemble model output statistics and minimum crps estimation. Monthly Weather Review 133(5): 1098–1118

Minor comment 1: Lack of technical information

All the results and conclusions from the study are dependent on the RF emulator. As such, I find that more information on the RF development and evaluation are needed. I highlight some aspects to prioritize here below.

- (a) The evaluation of the RF (L174-183) is assessed through a random sampling evaluation, but I find the details about the evaluation method somewhat unclear. First, the authors mention the "iteration of the procedure" (L180). However, it is not explained what is iterative in this procedure. Later in the manuscript, the authors often refer to "25 validation tests" (e.g., caption of Figure 4). But this number of 25 is not explained in the description of the evaluation method. Thus, I can only assume that the random validation is iterated 25 times. Second, it is unclear what the validation performance measure shown in Figure 4 represents. In Figure 4a, there are clearly much more than 25 points, but clearly less than $25 \times 55 = 1375$ points (where 55 is the number of test samples mentioned on L180). Thus, what does each point represent? In addition, why are there much less points shown in Figure 4b than Figure 4a? Finally, the authors explain that there are 55 test samples, but they draw 5 samples for 10% ranges between 0 and 100% (L179-180). As such, there should be $5 \times 10 = 50$ test samples I believe, not 55.
- (b) I wonder why it was decided to use this random evaluation procedure. In particular, the commonly-used 10-fold cross validation procedure would have been a more natural choice. This would also avoid the influence of sampling biases related with the random sampling of relatively few experiments (55 from the 1303 experiments per iteration). Since 10-fold cross validation was used for parameter fitting (L197), I suppose that there is no computational issue for this. Also, it would be straightforward to exclude the members from the 9 training folds as required by the specific experiments (e.g., exclude all SSP5-8.5 when training for woSSP585). Thus, is there any reason to prefer the random evaluation over the 10-fold cross validation? We reply here to comments (a) and (b). Referee #2 is right to say that 10-fold cross validation would be a more natural choice. The reason for proposing an alternative validation procedure is to make sure to reflect the ability of the emulator to perform well over a wide range of GSAT values instead of randomly selected cases. This ability is important in our case, because we discuss the performance with respect to the probabilistic projections given fixed GSAT values. The new Figure 5 (Major comment 2) illustrates this type of analysis.



New Figure 5: Boxplot of the RAE (a), Q^2 (b) and CRPS (c) performance indicator for different ranges of GSAT (indicated on the x-axis). The lower RAE and the closer Q^2 to one, the higher the emulator predictive capability. The

lower *CRPS*, the higher quality of the emulator predictive probabilistic distribution. The horizontal red dashed line indicates the median value calculated over all validation tests defined through the repeated validation procedure described in Sect. 2.4.1 considering the whole range of GSAT.

Though our GSAT definition does not strictly correspond to the global warming level defined in AR6, they can help end users to interpret the projections associated to temperature constraints as illustrated by recent projections for France by Le Cozannet et al. (2025).

We also thank Referee #2 for noticing the problem with the number of test cases. The presentation in Sect. 2.4.1 (on page 10, lines 201-211 has been clarified as follows: "In this study, we are more particularly interested in the ability of the emulator to perform well over a wide range of GSAT values. This is important in our case, because constraining the predictions to temperature constraints can help end-users to interpret the projections as illustrated by recent projections for France by Le Cozannet et al. (2025), although it should be noted that our GSAT definition does not strictly correspond to the global warming level (GWL) defined in AR6. Therefore, instead of relying on the widely used cross validation procedure (Hastie et al., 2009), we propose an alternative validation procedure adapted to our objective as follows: (1) the GSATs are classified into a finite number of intervals, the ends of which are defined by the GSAT percentiles, with levels ranging from 0 to 100% with a fixed increase of 25%. This results in the following GSAT intervals, $[0.705, 2.14^{\circ}C]$, $[2.14, 3.34^{\circ}C]$, $[3.34, 3.83^{\circ}C]$, and $[3.83, 5.00^{\circ}C]$; (2) for each interval, 50 samples are randomly selected. For one iteration of the procedure, a total of n_{test} =200 test samples are randomly selected. The procedure is repeated 25 times."

Added reference

Goneri Le Cozannet, Remi Thieblemont, Jeremy Rohmer and Cecile Capderrey (2025). Sealevel scenarios aligned with the 3rd adaptation plan in France. (in press) https://doi.org/10.5802/crgeos.290

(c) More technical details about the RF emulator construction would be beneficial. In particular, mixing categorical and continuous inputs is not straightforward, and may incur performance sensitivity to the RF design. For example, what is the splitting criterion used: mean absolute error, mean squared error, other? And how did the authors alleviate the potential issue of selection bias towards the inputs that have more possible splits? This could partly influence the different sensitivities to, for example, SSP5-8.5 scenario (global annual mean surface air temperature change, GSAT, is a continuous input with many different values), ISM (categorical input), κ (continuous input with few different values). As such, some information on these technical aspects would help the reader understand how modeling challenges may affect the results or not.

We agree that more technical details should be added.

We use the mean squared error in the loss function of the random forest model. The treatment of categorical variables is based on the recommendation by Hastie et al. (2009): chapter 9.2.4. We follow the implementation proposed by Wright et al. (2019), who showed that ordering the factor levels a priori, here by their mean response, is at least as good as the standard approach of considering all 2-partitions in all datasets considered, while being computationally faster; It has been shown to be more efficient than dummy coding and simply ignoring the nominal nature of the predictors as well.

However, as Referee #1 highlighted in her/his first comment, "the most readers of this journal are likely geoscientists, who will primarily be interested in the study's results". We propose to move these details in Appendix A on the Random Forest implementation.

Added references

Hastie, T., Tibshirani, R., Friedman, J. (2009). The Elements of Statistical Learning. Springer, New York. 2nd edition.

Wright, M. N., & König, I. R. (2019). Splitting on categorical predictors in random forests. PeerJ, 7, e6339.

Minor comment 2: Use of global mean temperature change

The authors aggregate all the combinations of emission scenario (SSP) and global climate model (GCM) as a value of GSAT. I wonder if this does not risk misrepresenting the climate forcing affecting the Greenland ice sheet (GrIS). In particular, a given GSAT could very well lead to different magnitudes of:

- (1) GrIS surface air temperature change
- (2) GrIS precipitation
- (3) GrIS ocean forcing

I expect that there may well be some substantial differences in these 3 components between different GCMs. It would be interesting to explore whether separating the single GSAT variable into these 3 separate components refines the emulator predictions.

This choice is based on the approach followed by previous emulation studies (Edwards et al. (2021). An alternative would use regional climate variables. Although this would improve the signal-to-noise ratio for the emulator, but would restrict us to using computationally expensive general circulation models from CMIP5/6, for which there only a few tens of models. With the GSAT option, as performed by Edwards et al. (2021), the simple climate model like FaIR can be also used to explore uncertainties in each scenario thoroughly, using the latest assessments of equilibrium climate sensitivity.

Referee #2 is also invited to refer to our reply to major comment 1 on the clarification of κ which is closely related to the sensitivity of the ocean forcing as a whole.

Note that we slightly change the definition of GSAT, in order to be more consistent with IPCC practices, by computing the difference between the temperature at the considered year and the mean temperature over the period 1995-2014.

Minor comment 3: Interpretation of some results

I find that the interpretation of results are not always well supported quantitatively. I note that, in some cases, this may simply be due to a lack of clarity in the interpretation. I provide here a few examples.

2.1 The Dh, DS definition

In Figure 6, the authors show the different combinations of decrease in MME size (DS) and deviations from original histograms (Dh) resulting from their model experiments. Firstly, I think that the manuscript would benefit from a clearer definition of Dh. It is defined as "the average difference in the count numbers between the two histograms (normalised by the total number of members)" (L172-173). I believe that the normalization is by the histogram counts, not the total number of members, because otherwise Dh would be proportional to DS. For example, assume that for a given variable, we have a hypothetical 3-category histogram with

counts 5, 10, 85 (i.e., n=100). In hypothetical experiment 1, the counts are 0, 10, 85 (i.e., n=95). In this case, DS = 100-5/100 = 0.95 and, following the definition, $Dh = 5+0+0/3 \times 1/100 = 1/60$. In hypothetical experiment 2, the counts are 5, 10, 80 (i.e., n=95). In this case, DS = 100-5/100 = 0.95 and $Dh = 0+0+5/3 \times 1/100 = 1/60$. This shows that taking "the average difference in the count numbers between the two histograms (normalised by the total number of members)" results in an identical pair (DS,Dh) for these two hypothetical experiments. I am probably misunderstanding here, but I think that a more precise definition would help. We thank Referee #2 for this insightful comment which made us rethink our procedure. As rightly shown by Referee #2, the proposed indicator D_h might fail to reflect the changes in the histograms.

Therefore, we propose to remove this analysis from the main text. In addition, we propose to support the discussion in Sect 4.1, on page 18, lines 352-359 with a complementary analysis (Supplementary materials S3) based on a well-established, and more widely used, criterion for comparing different probability distributions, namely the Kolomogorov-Smirnov (KS) criterion, instead of a criterion constructed from 'scratch'.

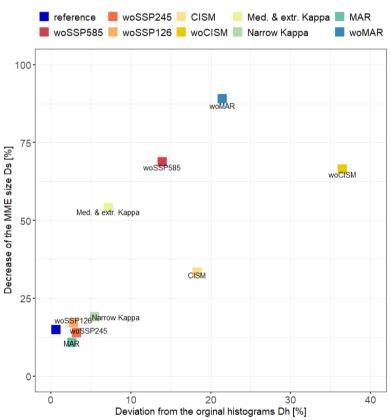


Figure S7: Position of the emulator's experiment in a (D_h, D_s) diagram where D_s measures the relative decrease in the MME size after applying the experiment, and D_h measures the deviation of the histograms from the original ones (see Supplementary Material S3). The blue-coloured marker refers to the reference solution defined as the mean value over the 25 iterations of the random validation exercise, described in Sect. 2.4, applied to the original dataset.

2.2 The Dh, DS results

I do not understand the interpretation of the impact from Dh, DS on the emulator performance (Sect. 3.3). First, the authors write "Excluding the extreme SSP scenario SSP5-8.5 (experiment 'woSSP585') has the largest impact in terms of RAE relative difference with respect to the original RF performance (Sect. 3.1), where RAE is increased of $\sim 10\%$ compared to the

original RAE value (Fig. 4)" (L245). However, Figure 7 shows a ~ 275% relative difference in RAE, so it is not clear to me where the value "~ 10%" comes from. Second, I do not follow the logic of the arguments. The authors write that (i) the high DS of woSSP585 causes large errors. But then, (ii) they argue that "this 'size effect' is not the only contributor to the performance impact, as shown by the 'woCISM' experiment, which removes an equivalent number of members to the 'woSSP585' experiment (Fig. 6), and the resulting RAE increase reaches half that of 'woSSP585' experiment" (L253). And (iii) that the woCISM experiment has the largest Dh value. However, when I interpret Figures 6 and 7, I find that (a) woSSP585 and woCISM have similar DS values (i.e., (ii)), (b) woCISM has higher Dh than woSSP585 (i.e., (iii)), but (c) that the errors from woSSP585 ar much higher than those of woCISM (Figure 7). So, it seems that the lower Dh of woSSP585 is accompanied by larger errors. This is the opposite message to that conveyed in the text: "This shows that the second important factor here is the diversity among the members within the MME after applying the experiment. The Dh indicator remains, however, a first-order approximation of this diversity (...)" (L256). The statement of greater diversity leading to lower errors, is not supported by the larger errors of woSSP585 compared to woCISM. To summarize: DS(woSSP585) ≈ DS(woCISM), Dh(woSSP585) < Dh(woCISM) where low Dh implies greater "diversity", but RAE(woSSP585) >> RAE(woCISM).

After reanalysing the results, we believe that the analysis Ds,Dh only supports the discussion in Sect. 4 but further developments may be needed to derive a robust Dh indicator. This is now clearly indicated on page 18, lines 352-359.

2.3 Figure S3 (in Section S2)

The authors write "The analysis of an alternative indicator of emulator's predictive capability in Supplementary materials S2 confirms these results" (L261). However, in my view, Figure 6 (RAE results) and Figure S3 (Q2, coefficient of determination results) show contrasting conclusions. For example, RAE of woCISM, CISM, and MAR are comparable (Figure 6). However, Q2 is clearly lower for woCISM than for MAR and CISM (Figure S3). This indicates differences when evaluating relative errors versus explained variance. Thus, these differences are potentially interesting to analyze, instead of being discarded as is done in the main text. In particular, they could relate to the emulator performance sensitivity to high versus low slc (the latter being more influential on relative metrics), or its sensitivity in the ability to predict values away from the mean value, or other aspects that would require investigation. Note that this links back to my general comment about the importance of understanding the RF emulator, because the interpretation of the results depends strongly on this understanding.

We thank Referee #2 for this valuable suggestion. We totally agree with Referee #2. We have included in the analysis not only RAE but also Q^2 . In addition we also propose to analyse the performance indicator \underline{CRPS} that measures the quality of the emulator's predictive probability as well. Referee #2 is invited also to refer to our detailed reply to Major comment 2.

2.4 Figure 8

There are many aspects that I find puzzling or questionable in Figure 8. Firstly, the results do not correspond to what is shown in Figure S4, where the Q5% and Q95% are shown with the black error bars. For example, in the column $\Delta GSAT=+3^{\circ}$, Q95% of woCISM, woSSP245, and woSSP585 are clearly strongly different from the Q95% labeled "original" (Figure S4). But Figure 8 shows that these differences are $\leq 1\%$. I believe that there is an inconsistency here, or something that I misunderstand about Figure 8.

Secondly, I do no understand how it is possible that the changes in median and quantiles at $\Delta GSAT = +4^{\circ}$ are so small for woSSP585. In this design experiment, the RF model has

presumably not even seen such levels of warming during training because the SSP 5-8.5 scenario has been excluded. But, by definition, tree models (including RF) predict slc based on decision rules seen during training. Thus, it is not clear how the RF can predict relatively similar slc values under $\Delta GSAT=+4\circ$ when excluding SSP 5-8.5 as when it is not excluded. I am probably misunderstanding something here, but I believe that the authors should explain this counter-intuitive aspect of their results.

We thank Referee #2 for this careful analysis. We confirm that, for some experiments, there are some discrepancies that have revealed a bug in our scripts for the plotting. We have updated the results in the new manuscript.

Minor comment 4: Some conclusions need to be put into perspective

For different aspects, I find that better communication and/or more context about the conclusions is needed. I highlight some key examples here.

(a) Concerning κ , the authors argue for "the lesser importance of the choice in the range of the Greenland tidewater glacier retreat parameter" (L21). However, they compare it with the influence of the SSP scenario and of the ISM choice. It is expected that a single parameter should have much less influence than a global warming scenario and than a full ice sheet model. We thank Referee #2 for this comment. Referee #2 is invited to refer the clarification made above about κ : It is not thought as a parameter of the ice-flow model, it rather represents the sensitivity of the ocean forcing as a whole. It may be thought of as defining the sensitivity of the downscaling from global model to local ice sheet scale, similar to the combined parameter choices in RCMs for downscaling climate conditions. In the studied MME, we have different RCMs, which have different sensitivities and produce different melt for the same global forcing. Since we have only one approach to 'downscale' the ocean forcing, κ is sampling that uncertainty in a similar way.

In addition, our previous study (Rohmer et al., 2022) highlighted the high importance of κ compared to other uncertainties. Referee #2 is also invited to refer to our detailed reply to major comment 1 as well to Referee #1's comment 3.

Reference

Rohmer, J., Thieblemont, R., Le Cozannet, G., Goelzer, H., and Durand, G.: Improving interpretation of sea-level projections through a machine-learning-based local explanation approach, Cryosphere, 16, 4637–4657, https://doi.org/10.5194/tc-16-4637-2022, 2022.

(b) It should be better emphasized that the probabilistic ranges shown by the authors are not probabilistic projections of Greenland slc. Instead, they show a range of emulator predictions (thus conditioned on the emulator architecture) assuming a uniform distribution over the different inputs (L186). Thus, it does not represent calibrated uncertainty accounting for model-observation misfits (e.g., Aschwanden and Brinkerhoff, 2022). And neither does it represent the slc PDF from the MME, because the uniform distribution over the input space is not representative of the MME itself (e.g., the minimum spatial resolution is clearly not uniform between 1 and 40 km, see Fig. 3). As such, I believe that the true meaning of the PDFs shown in Figure 5 should be explained explicitly in order to avoid any reader misinterpreting those PDFs.

We agree with this comment. To do so, we propose to clarify the caption of the new Fig. 5 as well as the description in the new Sect. 2.4.2 on page 11, lines 227-230.

(c) The authors make a conclusion on "the utmost importance of including the SSP5-8.5 scenario, due to the large number of simulations available and the range of global warming they cover" (L19-20). However, I do not think that the authors have proven the co-existence of these two points. For example, could it be that including only a few training simulations with high global warming forcing would be sufficient to drastically decrease the errors of woSSP585 shown in Figure 7? In other words, maybe the emulator needs only a few high-warming training examples to correctly interpolate in the existing range of warming scenarios. Or maybe, as the authors write (L19-20), it is also the high number of experiments that is important. However, as far as I understand, the results presented in this study do not allow to evaluate the relative importance of these two aspects.

We agree with Referee #2: we believe that the influence of the MME size is shown by our results, but disentangling this effect from the range of global warning remains too complicated at least with the procedure proposed here.

More broadly, we have nuanced our conclusions in this sense (Sect. 5) by outlining that our results depend on the considered MME. The title has also been modified in this sense, i.e., "Lessons for multi-model ensemble design drawn from emulator experiments: application to a large ensemble for future sea level contributions of the Greenland ice sheet."

References

Andy Aschwanden and DJ Brinkerhoff. Calibrated mass loss predictions for the greenland ice sheet. Geophysical Research Letters, 49(19):e2022GL099058, 2022.

Nicolai Meinshausen and Greg Ridgeway. Quantile regression forests. Journal of machine learning research, 7(6), 2006.

Lucas Mentch and Giles Hooker. Quantifying uncertainty in random forests via confidence intervals and hypothesis tests. Journal of Machine Learning Research, 17(26):1–41, 2016.

Christopher KI Williams and Carl Edward Rasmussen. Gaussian processes for machine learning, volume 2. MIT press Cambridge, MA, 2006.

We thank Referee #2 for the suggested references.

Orleans, July 4th, 2025 J. Rohmer¹ on behalf of the co-authors

¹ BRGM, 3 av. C. Guillemin - 45060 Orléans Cedex 2 – France