

Supplement of Geosci. Model Dev., 7, 1933–1943, 2014  
<http://www.geosci-model-dev.net/7/1933/2014/>  
doi:10.5194/gmd-7-1933-2014-supplement  
© Author(s) 2014. CC Attribution 3.0 License.



*Supplement of*

**Probabilistic calibration of a Greenland Ice Sheet model using spatially resolved synthetic observations: toward projections of ice mass loss with uncertainties**

**W. Chang et al.**

*Correspondence to:* W. Chang (wonchang@psu.edu)

# 1. Gaussian process emulator for principal components

In this section, we outline our statistical approach for ice sheet model emulation using Gaussian process (GP) models and principal component (PC) analysis (often referred to as empirical orthogonal functions, EOFs). Our approach follows Chang et al. (2014) in that we summarize the ice sheet model runs as PCs and calibrate the ice sheet parameters based on GP emulators for PCs. Our description of methods below therefore also closely follows the notation and description in Chang et al. (2014). By decomposing spatial patterns into a small number of variables representing important characteristics of model runs, our approach drastically increases computational efficiency without causing significant information loss.

We denote the number of model runs by  $p$  and the number of spatial locations spatial locations by  $n$ . For the SICOPOLIS model output (from Applegate et al. 2012) we use here,  $p = 99$  and  $n = 264$ . Note that the original ensemble in (Applegate et al. 2012) contains 100 model runs, but we leave one model run out to construct the synthetic truth. We let  $Y(\boldsymbol{\theta}, \mathbf{s})$  denote the ice thickness from the ice sheet model at a parameter setting  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_5)^T$  and a spatial location  $\mathbf{s}$ . We let  $\mathbf{s}_1, \dots, \mathbf{s}_n$  be the spatial locations of the model grid points and  $\mathbf{Y}(\boldsymbol{\theta}) = (Y(\boldsymbol{\theta}, \mathbf{s}_1), \dots, Y(\boldsymbol{\theta}, \mathbf{s}_n))$  be the vector of model output at a parameter setting  $\boldsymbol{\theta}$ . Let  $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p$  be the vectors of input parameters for our model.  $\mathbf{Y}$  is an  $n \times p$  matrix of the ice sheet model output where its rows correspond to spatial locations and columns to parameter settings, i.e.

$$\mathbf{Y} = \begin{pmatrix} Y(\boldsymbol{\theta}_1, \mathbf{s}_1), & Y(\boldsymbol{\theta}_2, \mathbf{s}_1), & \dots, & Y(\boldsymbol{\theta}_p, \mathbf{s}_1) \\ Y(\boldsymbol{\theta}_1, \mathbf{s}_2), & Y(\boldsymbol{\theta}_2, \mathbf{s}_2), & \dots, & Y(\boldsymbol{\theta}_p, \mathbf{s}_2) \\ \vdots, & \vdots, & \ddots, & \vdots \\ Y(\boldsymbol{\theta}_1, \mathbf{s}_n), & Y(\boldsymbol{\theta}_2, \mathbf{s}_n), & \dots, & Y(\boldsymbol{\theta}_p, \mathbf{s}_n) \end{pmatrix}.$$

Similarly,  $Z(\mathbf{s})$  denotes the observed ice sheet thickness at a location  $\mathbf{s}$ , and  $\mathbf{Z} = (Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n))^T$  is the  $n \times 1$  vector of the observational data. For the matrices and the statistical parameters used in the following sections, the subscript  $y$  indicates that a symbol is used for the emulation model, while the subscript  $d$  shows that a symbol is for the discrepancy model.

## 24 **2. Principal component analysis for model output**

25 The first step is summarizing the model output by principal component analysis. Using  
26 principal components in emulation and calibration has the following advantages: First, the  
27 principal components are independent of each other under Gaussian assumption, and this  
28 independence allows a substantial gain in computational efficiency because it enables us  
29 to treat each of the principal components separately in the emulation stage. Second, the  
30 principal components are the “best” summary of the model runs in the sense that they show  
31 the clearest contrast between model runs, among all possible linear combinations of model  
32 runs (see below for more details).

33 Following the standard procedure of principal component analysis, the column means are  
34 subtracted from each element in the corresponding columns such that each column is centered  
35 on zero. We apply singular value decomposition to this centered output matrix to find the  
36 scaled principal basis vectors  $\mathbf{k}_1 = \sqrt{\lambda_1}\mathbf{e}_1, \dots, \mathbf{k}_p = \sqrt{\lambda_p}\mathbf{e}_p$ , where  $\lambda_1 > \lambda_2 > \dots > \lambda_p$   
37 and  $\mathbf{e}_1, \dots, \mathbf{e}_p$  are ordered eigenvalues and their eigenvectors respectively. Each eigenvalue  
38 represents the explained variation for the corresponding principal component. We keep only  
39 the first  $J \ll p$  PCs with the largest explained variation (i.e. the largest eigenvalues) to  
40 minimize the information loss due to dimension reduction. The principal components for  
41 model output can be computed by

$$\mathbf{Y}^R = (\mathbf{K}_y^T \mathbf{K}_y)^{-1} \mathbf{K}_y^T \mathbf{Y} = (\mathbf{Y}_1^R \dots \mathbf{Y}_J^R)^T$$

42 where  $\mathbf{K}_y = (\mathbf{k}_1, \dots, \mathbf{k}_J)$  is the principal basis matrix.  $\mathbf{Y}_i^R = (Y_i^R(\boldsymbol{\theta}_1), \dots, Y_i^R(\boldsymbol{\theta}_p))^T$  is the  
43  $p \times 1$  vector of the  $i$ th principal components, and  $Y_i^R(\boldsymbol{\theta}_j)$  is the  $i$ th principal component at  
44 the parameter setting  $\boldsymbol{\theta}_j$ . The resulting matrix  $\mathbf{Y}^R$  is the summarized output matrix with  
45 rows for PCs and columns for parameter settings. The procedure reduces the size of the  
46 data from  $n \times p$  to  $J \times p$ .

### 47 3. Gaussian process emulator

48 We emulate the ice sheet model output using Gaussian processes (GP), a fast method  
49 for probabilistic interpolation between existing model runs (Sacks et al. 1989; Kennedy and  
50 O’Hagan 2001; Higdon et al. 2008; Drignei et al. 2008; Rougier 2008; Goldstein and Rougier  
51 2009; Bhat et al. 2012; Holden et al. 2010; Lee et al. 2011; Olson et al. 2012, 2013; Mc-  
52 Neall et al. 2013; Williamson et al. 2013). The GP emulator approach yields a flexible  
53 approximation without requiring detailed physical information on the ice sheet model, un-  
54 like linear regression-based emulators (cf. Piani et al. 2005). By interpolating existing model  
55 runs at different parameter settings, a GP emulator provides a reasonable approximation to  
56 the original model unless the model output abruptly changes in the input parameter space.  
57 Interpolation using GP emulator is essentially kriging in the input parameter space; the  
58 interpolator is a random process with (i) a mean term that is the optimal interpolation be-  
59 tween ice sheet model runs in terms of the expected mean squared error and (ii) a variance  
60 term that quantifies the uncertainty of the interpolation. In this section we first describe the  
61 basics of the GP emulator approach and explain how to construct an emulator for principal  
62 components.

63 Before describing emulation for principal components, we illustrate the GP emulation  
64 approach with  $n = 1$ , where the ice model output at each parameter setting is a scalar. For  
65 ease of exposition we further simplify the case by assuming that the input parameter is also  
66 a scalar. The collection of the model output is simply a vector  $\mathbf{Y} = (Y(\theta_1), \dots, Y(\theta_p))^T$  at  
67 the parameter settings  $\theta_1, \dots, \theta_p$ . The GP emulator for the model output is then given by

$$\mathbf{Y} \sim N(\mathbf{0}, \Sigma),$$

68 with a  $p \times p$  covariance matrix  $\Sigma$ . While linear regression finds a trend that fits the data  
69 well using a mean function, the GP model interpolates the data using a covariance structure.  
70 Unlike the linear regression, which requires careful specification of the mean function along  
71 with various statistical assumptions when dealing with highly nonlinear processes such as ice

72 model outputs, the GP model can handle such processes using a relatively simple covariance  
 73 function. Only required assumption for the GP model is that the model output is a smoothly  
 74 varying curve in the parameter space without too many abrupt changes. A common choice  
 75 for the covariance function for a GP emulator is the squared exponential covariance, which  
 76 defines the  $(i, j)$ th element of  $\Sigma$  as

$$\text{Cov}(Y(\theta_i), Y(\theta_j) | \zeta, \kappa, \phi) = \zeta 1(\theta_i = \theta_j) + \kappa \exp\left(-\left(\frac{\theta_i - \theta_j}{\phi}\right)^2\right),$$

77 with all positive  $\zeta$ ,  $\kappa$  and  $\phi$ . The range parameter  $\phi$  defines how fast the model output  
 78 is changing as the value of the parameter changes. The partial sill specifies the overall  
 79 magnitude of the process, and the nugget parameter  $\zeta$  captures the variability caused by  
 80 various sources other than input parameters, such as the effect of initial conditions.

81 We now turn our attention to emulation for the principal components. Because the prin-  
 82 cipal components are uncorrelated with each other by construction, we can model each of  
 83 them separately using independent GPs. Note that this procedure ignores the dependence  
 84 between the principal components that is not captured by the covariances. However, ac-  
 85 cording to our cross-validation experiments for various models including SICOPOLIS, the  
 86 emulator based on this assumption usually provides an accurate approximation to the origi-  
 87 nal model (see e.g. Chang et al. 2014, Figure 2). We model each  $\mathbf{Y}_i^R$  using a GP with mean  
 88 zero and covariance determined by the following squared exponential covariance function:

$$\text{Cov}(Y_i^R(\boldsymbol{\theta}_j), Y_i^R(\boldsymbol{\theta}_k); \zeta_i, \kappa_{y,i}, \phi_i) = \zeta_i 1(\boldsymbol{\theta}_j = \boldsymbol{\theta}_k) + \kappa_{y,i} \exp\left(-\sum_{l=1}^5 \left(\frac{\theta_{jl} - \theta_{kl}}{\phi_{il}}\right)^2\right),$$

89 where  $\zeta_i, \kappa_{y,i}, \phi_{i1}, \dots, \phi_{i5} > 0$  are covariance parameters,  $\theta_{jl}$  is the  $l$ th element of  $\boldsymbol{\theta}_j$ , and  $1(\cdot)$  is  
 90 the index function. The covariance parameters  $(\zeta_1, \kappa_{y,1}, \phi_{11}, \dots, \phi_{15}), \dots, (\zeta_J, \kappa_{y,J}, \phi_{J1}, \dots, \phi_{J5})$   
 91 are estimated by maximum likelihood estimation (MLE). Our emulator, denoted by  $J \times 1$   
 92 vector-valued function  $\boldsymbol{\eta}(\boldsymbol{\theta}, \mathbf{Y}^R)$ , is the predictive distribution of PCs at an untried param-  
 93 eter setting  $\boldsymbol{\theta}$  defined by the fitted GPs. Using the PC emulator, we can also emulate the  
 94 original model transect by computing  $\mathbf{K}_y \boldsymbol{\eta}(\boldsymbol{\theta}, \mathbf{Y}^R)$ .

95 Note that our approach allows significant improvements in computational efficiency.  
 96 Without any dimension reduction, the computational cost for a single likelihood evaluation  
 97 scales as  $\mathcal{O}(n^3p^3)$ , which corresponds to a few hours of computing time. Thus, application  
 98 of any numerical methods requiring repeated evaluation of the likelihood function is com-  
 99 putationally prohibitive if no dimensional reduction is performed. Our approach decreases  
 100 the computational complexity to  $\mathcal{O}(Jp^3)$ , and this is a reduction from  $3.18 \times 10^{14}$  flops to  
 101  $1.56 \times 10^8$  flops when using 10 principal components ( $J = 10$ ). Using 10 principal compo-  
 102 nents captures more than 90% of the variation in the model output, and we have confirmed  
 103 that using more than 10 principal components does not significantly improve the emulation  
 104 accuracy by cross-validation.

## 105 4. Model parameter calibration

106 In this section, we formulate the probability model for calibration using the PC emula-  
 107 tor constructed above and explain the inference procedure for the model parameters using  
 108 Markov chain Monte Carlo (MCMC). The main goal here is to estimate the input param-  
 109 eters using a probability model that combines the model for the emulator described above  
 110 and a discrepancy term that detects the systematic model-observation discrepancy. The  
 111 calibration model operates in a reduced dimensional space defined by the principal compo-  
 112 nents computed above and the kernel convolution (Higdon et al. 2008, explained below) and  
 113 allows us to avoid computational issues for dealing with high-dimensional data (Chang et al.  
 114 2014). Based on this model, one can efficiently sample from the posterior density of input  
 115 parameters by MCMC.

116 We assume that the observational dataset is emulator output contaminated by model  
 117 discrepancy and observational error;

$$\mathbf{Z} = \mathbf{K}_y \boldsymbol{\eta}(\boldsymbol{\theta}^*, \mathbf{Y}^R) + \mathbf{K}_d \boldsymbol{\nu} + \boldsymbol{\epsilon}, \quad (\text{S1})$$

118 where  $\boldsymbol{\theta}^*$  is the best fit input parameter setting (Bayarri et al. 2007; Rougier 2007) for the

119 observational data, and  $\boldsymbol{\epsilon} \sim N(0, \sigma^2 \mathbf{I}_n)$  is the observational error with variance  $\sigma^2 > 0$ .  
 120  $\mathbf{K}_d \boldsymbol{\nu}$  is the model-observation discrepancy picking up systematic differences between the  
 121 model and the observations (cf. Bayarri et al. 2007; Bhat et al. 2012), where  $\mathbf{K}_d$  is a kernel  
 122 basis matrix relating the spatial locations  $\mathbf{s}_1, \dots, \mathbf{s}_n$  to  $J_d$  knot locations  $\mathbf{a}_1, \dots, \mathbf{a}_{J_d}$ , and  
 123  $\boldsymbol{\nu} \sim N(\mathbf{0}, \kappa_d \mathbf{I}_{J_d})$  is the vector of knot processes, a set of random variables assigned to each of  
 124 the knot locations with variance  $\kappa_d > 0$ . Our choice for the kernel function is an exponential  
 125 covariance given by

$$\{\mathbf{K}_d\}_{ij} = \exp\left(-\frac{|\mathbf{s}_i - \mathbf{a}_j|}{\phi_d}\right),$$

126 with  $\phi_d > 0$ . The variance parameter  $\kappa_d$  is subject to inference, and the correlation pa-  
 127 rameter  $\phi_d$  is pre-specified by expert judgment. In our implementation, we choose  $\phi_d$  as  
 128 5% of the maximum distance between the spatial locations on the model grid to yields a  
 129 sufficiently flexible discrepancy pattern. Fixing the range parameter not only reduces the  
 130 computational cost for likelihood computation but also improves the identifiability between  
 131 the input parameters and the discrepancy process. Note that the kernel basis often needs  
 132 to be substituted by its scaled principal basis (eigenvectors) to improve identifiability; see  
 133 Chang et al. (2014) for a more detailed discussion. We used the 30 leading components  
 134 for  $\mathbf{K}_d$  in our implementation. We apply a similar dimension reduction described in the  
 135 previous section to find  $\mathbf{Z}^R$ , a summary of the observed transect as follows:

$$\mathbf{Z}^R = (\mathbf{K}^T \mathbf{K})^{-1} \mathbf{K}^T \mathbf{Z}, \quad (\text{S2})$$

136 and therefore the model for  $\mathbf{Z}^R$  can be written as

$$\mathbf{Z}^R \sim N\left(\begin{pmatrix} \boldsymbol{\mu}_\eta \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \Sigma_\eta & \mathbf{0} \\ \mathbf{0} & \kappa_d \mathbf{I}_{J_d} \end{pmatrix} + \sigma^2 (\mathbf{K}^T \mathbf{K})^{-1}\right),$$

137 where  $\boldsymbol{\mu}_\eta$  and  $\Sigma_\eta$  are the mean and covariance, respectively, of the emulator  $\boldsymbol{\eta}(\boldsymbol{\theta}^*, \mathbf{Y}^R)$ , and  
 138  $\mathbf{K} = (\mathbf{K}_y \ \mathbf{K}_d)$ .

139 The parameters to be estimated in the calibration model are the ice sheet model input  
 140 parameters  $\boldsymbol{\theta}^*$ , the discrepancy parameter  $\kappa_d$ , and the observational error variance  $\sigma^2$ . We

141 also re-estimate the partial sill parameters  $\boldsymbol{\kappa}_y = (\kappa_{y,1}, \dots, \kappa_{y,J})$  for the emulator (Bayarri  
 142 et al. 2007; Bhat et al. 2012; Chang et al. 2014). This allows the emulator process to be re-  
 143 scaled to better match the observational data. We define the posterior density based on the  
 144 likelihood function given by (S2) denoted by  $\ell(\mathbf{Z}^R|\boldsymbol{\theta}^*, \boldsymbol{\kappa}_y, \kappa_d, \sigma^2, \mathbf{Y}^R)$  and some standard prior  
 145 specifications denoted by  $f(\boldsymbol{\theta}^*)$ ,  $f(\boldsymbol{\kappa}_y)$ ,  $f(\kappa_d)$ , and  $f(\sigma^2)$  (Higdon et al. 2008; Chang et al.  
 146 2014). Each of the input parameters in  $\boldsymbol{\theta}^*$  receives a flat prior on a broad range determined  
 147 by model ensemble design and physical knowledge. The observational error variance  $\sigma^2$  and  
 148 the variance for the discrepancy  $\kappa_d$  have non-informative inverse-gamma priors with small  
 149 shape parameters. We specify somewhat informative priors for  $\kappa_{y,1}, \dots, \kappa_{y,J}$  by specifying  
 150 a large shape parameter in order to avoid numerical instability and identifiability issues  
 151 (Higdon et al. 2008). The posterior distribution resulting from the above model is

$$\pi(\boldsymbol{\theta}^*, \boldsymbol{\kappa}_y, \kappa_d, \sigma^2|\mathbf{Z}^R, \mathbf{Y}^R) \propto \ell(\mathbf{Z}^R|\boldsymbol{\theta}^*, \boldsymbol{\kappa}_y, \kappa_d, \sigma^2, \mathbf{Y}^R) f(\boldsymbol{\theta}^*) f(\boldsymbol{\kappa}_y) f(\kappa_d) f(\sigma^2),$$

152 where

$$\begin{aligned} \ell(\mathbf{Z}^R|\boldsymbol{\theta}^*, \boldsymbol{\kappa}_y, \kappa_d, \sigma^2, \mathbf{Y}^R) &\propto |\Sigma_{\boldsymbol{\eta}} + \mathbf{K}^T \mathbf{K} \sigma^2|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} \mathbf{Z}^{RT} (\Sigma_{\boldsymbol{\eta}} + \mathbf{K}^T \mathbf{K} \sigma^2)^{-1} \mathbf{Z}^R\right) \\ f(\boldsymbol{\theta}^*) &\propto 1(\boldsymbol{\theta}^* \in \Theta), \Theta \text{ represents the range of } \boldsymbol{\theta}, \\ f(\boldsymbol{\kappa}_y) &\propto \prod_{i=1}^J \kappa_{y,i}^{-a_{y,i}-1} \exp\left(-\frac{b_{y,i}}{\kappa_{y,i}}\right), a_{y,1}, \dots, a_{y,J}, b_{y,1}, \dots, b_{y,J} > 0 \\ f(\kappa_d) &\propto \kappa_d^{-a_d-1} \exp\left(-\frac{b_d}{\kappa_d}\right), a_d, b_d > 0 \\ f(\sigma^2) &\propto \sigma^{-2(a_{\sigma}+1)} \exp\left(-\frac{b_{\sigma}}{\sigma^2}\right), a_{\sigma}, b_{\sigma} > 0. \end{aligned}$$

153 For each  $i$ , we set  $a_{y,i} = 50$  and choose  $b_{y,i}$  such that the mode of the prior density  $b_{y,i}/(a_{y,i}+1)$   
 154 coincides with the MLE of  $\kappa_{y,i}$  computed in the emulation stage. For other parameters, we  
 155 impose vague priors by setting  $a_d = 2$ ,  $b_d = 3$ ,  $a_{\sigma} = 2$ , and  $b_{\sigma} = 3$ .

156 The synthetic observations used in our perfect model experiment are constructed by  
 157 superimposing a random error generated from a Gaussian process model on the assumed  
 158 true ice sheet status. To make our experiment more realistic, the discrepancy process is  
 159 generated from a different model to the discrepancy term that we use in the equation (S1).



160 The covariance function that we use for the Gaussian process model for simulated discrepancy  
161 here is a squared exponential covariance having range of 2100 km, partial sill of 2500 m, and  
162 a nugget of 1 m. Our choice for the simulated discrepancy process is based on the following  
163 two general assumptions: (i) the discrepancy is statistically identifiable from the emulator  
164 process, and (ii) SICOPOLIS has an enough skill to reproduce the observed ice profile. (i)  
165 is related to the value of the range parameter, which controls the effective distance at which  
166 two spatial locations are uncorrelated. To ensure that the discrepancy process is identifiable  
167 from the emulator process, we set the range parameter to be very large (80% of the spatial  
168 range of the model output) so that the discrepancy operates in a different spatial scale  
169 to the emulator process. (ii) is related to the value of the partial sill, which defines the  
170 magnitude of the discrepancy. Here we let the value of the partial sill to be reasonably small  
171 to simulate the situation in which the structural error is not large and therefore SICOPOLIS  
172 can reproduce the observed ice profile reasonably well. Note that calibration based on any  
173 framework, including our approach, may yield misleading results if any of the assumptions  
174 are violated. For example, if the discrepancy process operates over a similar spatial scale  
175 to the emulator process (i.e. (i) does not hold), the discrepancy causes identifiability issues  
176 and hence introduces a significant bias in the calibration result. If the magnitude of the  
177 discrepancy is too large (i.e. (ii) does not hold) compared to the variation between model  
178 outputs, the calibration results will become essentially non-informative (i.e., the likelihood  
179 then has little effect on the posterior distribution). Note that these are common issues for  
180 calibration methods in general. The curves in Figure S1 show that realizations from the  
181 discrepancy process are clearly identifiable from the difference between model runs.

182 Based on the pseudo-observations, we infer the parameters using the MCMC sample  
183 from the above posterior distribution obtained by the Metropolis-Hastings algorithm (cf.  
184 Higdon et al. 2009). In particular, we infer the input parameters in  $\boldsymbol{\theta}^*$  by investigating  
185 their marginal density  $\pi(\boldsymbol{\theta}^* | \mathbf{Z}^R, \mathbf{Y}^R)$ . In our perfect model experiment, we obtained 300,000  
186 draws using block updating when estimating the full joint density of all five parameters. The

187 computing time takes about eight hours on a single high-performance core. For inference on  
188 individual input parameter, only 300,000 draws using block updating is sufficient. In both  
189 cases, we confirmed that the Monte Carlo chain is well-mixed by comparing the densities of  
190 the first half of the chain with the entire chain. We find the probability density of the input  
191 parameters via kernel density estimation for the MCMC sample. The estimated density can  
192 be easily plotted for visual analysis as shown in Figures 3 and 4. Note that ignoring the  
193 spatially correlated discrepancy results in a notably biased calibration results in our perfect  
194 model experiment. See Figure S2 for a comparison of posterior densities with and without  
195 the discrepancy term.

## 196 **5. Ice volume change projection based on calibrated pa-** 197 **rameters**

198 One important purpose of parameter calibration is making better projections for the  
199 future ice sheet mass loss. In our illustrative example, the variable that we want to project  
200 is the ice volume change from present to 2100 in meters of sea level equivalent. For each  
201 model run, we compute the ice volume change by subtracting the current ice volume from the  
202 future ice volume. Making future projections based on calibration results requires a function  
203 that relates input parameter values  $\theta^*$  to future changes in ice sheet volume. We construct  
204 such a function by interpolating the ice volume changes for the 100 input parameter settings  
205 in our ensemble. Using this function we convert the MCMC chain for the input parameters  
206 generated in the previous section into Monte Carlo sample for the future ice volume changes.

207 Among many possible choices for the interpolator, we use the Gaussian process emulator  
208 similar to the model described in 3. More specifically, we fit a Gaussian process model for  
209 the ice volume change over the input parameter space with zero-mean and the covariance

210 function

$$Cov(\Delta v(\boldsymbol{\theta}_j), \Delta v(\boldsymbol{\theta}_k); \zeta^{vol}, \kappa^{vol}, \phi^{vol}) = \zeta^{vol} 1(\boldsymbol{\theta}_j = \boldsymbol{\theta}_k) + \kappa^{vol} \exp\left(-\sum_{l=1}^5 \frac{|\theta_{jl} - \theta_{kl}|}{\phi_l^{vol}}\right),$$

211 for any given design points  $\boldsymbol{\theta}_j$  and  $\boldsymbol{\theta}_k$  ( $j, k = 1, \dots, 100$ ), where  $\Delta v(\boldsymbol{\theta})$  is the volume change  
212 at a parameter setting  $\boldsymbol{\theta}$ , and  $\zeta^{vol}, \kappa^{vol}, \phi_1^{vol}, \dots, \phi_5^{vol} > 0$  are the covariance parameters  
213 that need to be estimated via MLE. The resulting function can predict ice volume change  
214 at any given value of  $\boldsymbol{\theta}$  as the conditional mean given by the standard kriging approach  
215 (Cressie 1993). Figure S3 shows the marginal surface of the projection as a function of input  
216 parameters. To validate the emulator constructed here, we have conducted leave-5%-out  
217 cross-validation. The mean error rate, computed by dividing the RMS by the overall mean,  
218 is around 16%; the error rate is a little higher than the heuristic upper limit for the generally  
219 acceptable emulation error (10%) due to the irregular behavior of the volume change surface.

220 We obtain a Monte Carlo sample of ice volume projections by supplying the posterior  
221 sample of the calibrated parameters to the interpolation function. Each element of the  
222 posterior sample is converted to ice volume change. The predictive density of the ice volume  
223 projection can be found by applying kernel density estimation. We find the prior density of  
224 the projections in the same manner; we convert the design points of the existing model runs  
225 into the ice volume changes and compute the predictive density for it using kernel density  
226 estimation.

## 227 6. Cross-validation

228 To investigate (i) whether the perfect model experiment results shown in the main text  
229 are sensitive to the values of input parameters assumed as the synthetic truth, and (ii)  
230 whether the prediction intervals for ice volume projections generated from our method have  
231 the right coverage, we have conducted leave-one-out cross-validation across all input param-  
232 eter settings in the ensemble. In other words, we have repeated the same perfect model  
233 experiment described in the previous sections for all 100 possible different synthetic truths.

234 We summarize the cross-validation results for emulation and calibrated projections in Figure  
235 S4 and Figure S5, respectively.

236 The results in Figure S4 show that our emulator can predict the model output reasonably  
237 well across all input parameter settings. The predicted ice volume thickness profiles are  
238 concentrated around the diagonal line that connects the lower left and the upper right corners  
239 of the plot, and hence the emulator can predict the model output reasonably well for most  
240 input parameter settings. Note that leave-one-out cross-validation is already rigorous enough  
241 in our case due to the sparsity of the design points (100 points in 5-dimensional space) for  
242 the input parameters in our ensemble. We have also conducted leave-10-out cross-validation  
243 for emulation and the results are essentially the same (not shown).

244 The plots in Figure S5 show that the prediction intervals generated from our approach  
245 achieve the nominal coverage level only when the modern ice volume generated by the syn-  
246 thetic truth is close enough to the observed volume (i.e. within 10% of the observed value).  
247 The width of the prediction interval also varies considerably across the different assumed  
248 truths. Therefore, consistent with the findings in McNeall et al. (2013), selection of the  
249 assumed truth affects the calibration performance. Another important observation is that  
250 including the discrepancy term reduces the overconfidence issue that occurs when the syn-  
251 thetic truths are outside of the 90-110% range. The prediction intervals are overconfident  
252 when the synthetic truth is outside of this range because the coverage is consistently less  
253 than 95%. Including the discrepancy term reduces this issue in some degree since it make the  
254 actual coverage closer to the nominal coverage when the synthetic truth yields the modern  
255 ice volume that is within at most 70% of the observed volume. However, this correction  
256 effect is not sufficient to make the prediction intervals achieve the nominal coverage.

257 The cross-validation results allow us to examine the interaction between input parameters  
258 across all possible choices of the synthetic truth. We have computed the rank correlations  
259 between the input parameters across all 100 ensemble members and summarized their distri-  
260 butions in Figure S6. From the shapes of the densities we can identify five pairs of parameters

261 that tend to be more negatively correlated: (i) the flow factor and the snow PDD factor, (ii)  
262 the flow factor and the geothermal heat flux, (iii) the basal sliding factor and the ice PDD  
263 factor, (iv) the geothermal heat flux and the ice PDD factor, and (v) the ice PDD factor  
264 and the snow PDD factor.

## 265 **7. Summary**

266 We describe an ice sheet model calibration approach based on PCs of the model output  
267 and the observational data. We build a GP emulator for the PCs of the model output as  
268 a fast approximation to the ice sheet model. The calibration model links the observed PCs  
269 with the input parameters using the GP emulator while taking the systematic discrepancy  
270 into account. We infer the input parameters along with other statistical parameters in  
271 the calibration model using MCMC. Combined with projections generated by the ice sheet  
272 model, the resulting posterior density of the parameters provide calibrated probabilistic  
273 projections of the future ice sheet volume changes. Our cross-validation results across all  
274 input parameter settings in the ensemble show that the probabilistic projections achieves  
275 the nominal coverage rate when the synthetic truth yields a modern ice volume that is close  
276 to the observed volume.

## REFERENCES

- 279 Applegate, P., N. Kirchner, E. Stone, K. Keller, and R. Greve, 2012: An assessment of  
280 key model parametric uncertainties in projections of Greenland Ice Sheet behavior. *The*  
281 *Cryosphere*, **6 (3)**, 589–606.
- 282 Bayarri, M., et al., 2007: Computer model validation with functional output. *The Annals of*  
283 *Statistics*, **35 (5)**, 1874–1906.
- 284 Bhat, K., M. Haran, R. Olson, and K. Keller, 2012: Inferring likelihoods and climate system  
285 characteristics from climate models and multiple tracers. *Environmetrics*, **23 (4)**, 345–362.
- 286 Chang, W., M. Haran, R. Olson, and K. Keller, 2014: Fast dimension-reduced climate  
287 model calibration and the effect of data aggregation. *The Annals of Applied Statistics*,  
288 **8 (2)**, 649–673.
- 289 Cressie, N., 1993: *Statistics for Spatial Data*. Wiley, New York.
- 290 Drignei, D., C. Forest, and D. Nychka, 2008: Parameter estimation for computationally  
291 intensive nonlinear regression with an application to climate modeling. *The Annals of*  
292 *Applied Statistics*, **2 (4)**, 1217–1230.
- 293 Goldstein, M. and J. Rougier, 2009: Reified bayesian modelling and inference for physical  
294 systems. *Journal of Statistical Planning and Inference*, **139 (3)**, 1221–1239.
- 295 Higdon, D., J. Gattiker, B. Williams, and M. Rightley, 2008: Computer model calibration us-  
296 ing high-dimensional output. *Journal of the American Statistical Association*, **103 (482)**,  
297 570–583.

- 298 Higdon, D., C. Reese, J. Moulton, J. Vrugt, and C. Fox, 2009: Posterior exploration for com-  
299 putationally intensive forward models. *Handbook of Markov Chain Monte Carlo*, S. Brooks,  
300 A. Gelman, G. Jones, and X. Meng, Eds., CRC Press, New York, 401–418.
- 301 Holden, P. B., N. Edwards, K. Oliver, T. Lenton, and R. Wilkinson, 2010: A probabilistic  
302 calibration of climate sensitivity and terrestrial carbon change in genie-1. *Climate dynam-*  
303 *ics*, **35 (5)**, 785–806.
- 304 Kennedy, M. and A. O’Hagan, 2001: Bayesian calibration of computer models. *Journal of*  
305 *the Royal Statistical Society. Series B (Statistical Methodology)*, **63 (3)**, 425–464.
- 306 Lee, L., K. Carslaw, K. Pringle, G. Mann, and D. Spracklen, 2011: Emulation of a com-  
307 plex global aerosol model to quantify sensitivity to uncertain parameters. *Atmospheric*  
308 *Chemistry and Physics*, **11 (23)**, 12 253–12 273.
- 309 McNeall, D., P. Challenor, J. Gattiker, and E. Stone, 2013: The potential of an observational  
310 data set for calibration of a computationally expensive computer model. *Geoscientific*  
311 *Model Development*, **6 (5)**, 1715–1728.
- 312 Olson, R., R. Sriver, M. Goes, N. Urban, H. Matthews, M. Haran, and K. Keller, 2012:  
313 A climate sensitivity estimate using Bayesian fusion of instrumental observations and an  
314 Earth System model. *Journal of Geophysical Research: Atmospheres*, **117 (D4)**, D04 103,  
315 doi:10.1029/2011JD016 620.
- 316 Olson, R., R. Sriver, M. Haran, W. Chang, N. Urban, and K. Keller, 2013: What is the  
317 effect of unresolved internal climate variability on climate sensitivity estimates? *Journal*  
318 *of Geophysical Research: Atmospheres*, **118 (10)**, 4348–4358.
- 319 Piani, C., D. Frame, D. Stainforth, and M. Allen, 2005: Constraints on climate change from  
320 a multi-thousand member ensemble of simulations. *Geophysical Research Letters*, **32 (23)**.

- 321 Rougier, J., 2007: Probabilistic inference for future climate using an ensemble of climate  
322 model evaluations. *Climatic Change*, **81** (3–4), 247–264.
- 323 Rougier, J., 2008: Efficient emulators for multivariate deterministic functions. *Journal of*  
324 *Computational and Graphical Statistics*, **17** (4), 827–843.
- 325 Sacks, J., W. Welch, T. Mitchell, and H. Wynn, 1989: Design and analysis of computer  
326 experiments. *Statistical Science*, **4** (4), 409–423.
- 327 Williamson, D., M. Goldstein, L. Allison, A. Blaker, P. Challenor, L. Jackson, and K. Ya-  
328 mazaki, 2013: History matching for exploring and reducing climate model parameter space  
329 using observations and a large perturbed physics ensemble. *Climate Dynamics*, **41** (7–8).



330 **List of Tables**

331 1 Summary of notation

17

TABLE 1. Summary of notation

Symbol	Definition
$p$	number of model runs
$n$	number of spatial locations for model grid
$\theta_{il}$	value of $l$ th input parameter in $\boldsymbol{\theta}_i$ for $i$ th model run
$\boldsymbol{\theta}_i$	input parameter setting for $i$ th model run
$\mathbf{s}_j$	$j$ th location on computer model grid
$Y(\boldsymbol{\theta}_i, \mathbf{s}_j)$	model output at location $\mathbf{s}_j$ for input parameter setting $\boldsymbol{\theta}_i$
$\mathbf{Y}$	$p \times n$ matrix of all model ensemble
$J$	number of principal components used in emulation
$\mathbf{Y}_i^R$	principal components for $i$ th model run, $(Y_i^R(\boldsymbol{\theta}_1), \dots, Y_i^R(\boldsymbol{\theta}_p))$
$\mathbf{Y}^R$	$p \times J$ matrix of all principal components, $(\mathbf{Y}_1^R, \dots, \mathbf{Y}_p^R)^T$
$\mathbf{k}_j$	$j$ th (scaled) principal component basis vector
$\mathbf{K}_y$	$n \times J$ principal component basis matrix $(\mathbf{k}_1, \dots, \mathbf{k}_J)$
$\boldsymbol{\eta}(\boldsymbol{\theta}, \mathbf{Y}^R)$	emulator for principal components
$\boldsymbol{\theta}^*$	true or fitted value of computer model parameter for observational data
$\zeta_i$	nugget for $i$ th principal component emulator
$\kappa_{y,i}$	partial sill for $i$ th principal component emulator
$\phi_{il}$	range for $i$ th principal component and $l$ th input parameter
$\mathbf{K}_d$	kernel basis matrix for discrepancy
$\mathbf{a}_j$	$j$ th knot location for kernel basis
$\boldsymbol{\nu}$	vector of knot processes
$J_d$	number of knot locations
$\kappa_d$	variance of knot processes in $\boldsymbol{\nu}$
$\phi_d$	range parameter for discrepancy kernel basis in $\mathbf{K}_d$
$\mathbf{Z}$	vector of observational data
$\mathbf{Z}^R$	principal components for observational data
$\boldsymbol{\epsilon}$	vector of observational errors
$\sigma^2$	variance of observational errors
$\mathbf{K}$	basis matrix for observational data, $\mathbf{K} = (\mathbf{K}_y \mathbf{K}_d)$
$\Delta v(\boldsymbol{\theta})$	volume change projection for input parameter setting $\boldsymbol{\theta}$
$\zeta^{vol}$	nugget for volume change emulator
$\kappa^{vol}$	partial sill for volume change emulator
$\phi_l^{vol}$	$l$ th range for volume change emulator

## 332 List of Figures

333 S1 Comparison between (i) residuals between the synthetic truth used in the  
334 main text (model run #67 in Applegate et al. 2012) and other model runs  
335 (black solid curves) and (ii) 30 different realizations from the model for the  
336 simulated discrepancy (red solid curves). The residuals are computed by sub-  
337 tracting the synthetic truth from each of the other model runs. For better  
338 display, we show only residual curves whose ranges are within (-500,500). It is  
339 easy to see that the black curves and red curves are generated from different  
340 processes, and therefore those two groups of curves can be separated by sta-  
341 tistical inference (hence the discrepancy is identifiable). The magnitudes of  
342 the simulated discrepancy processes are well within the range covered by the  
343 model runs (hence the posterior density of input parameters does not show  
344 too large variation). 21

345 S2 Comparison between calibration results with and without the discrepancy  
346 term  $\mathbf{K}_d \boldsymbol{\nu}$  in the calibration model in (S1). In each panel, we tried to learn  
347 each of the parameters while fixing the other parameters at their assumed-  
348 true values. The prior densities are assumed to be uniform over a broad range  
349 (dashed red lines). While the posterior densities computed by including the  
350 discrepancy term in the model (solid black curves) pick up the true parameter  
351 values without notable biases, the posterior densities without the discrepancy  
352 term (solid blue curves) cannot recover the true values. 22

353 S3 Surfaces of ice volume change projections between 2005 and 2100 projected  
354 onto marginal spaces of all pairs of input parameters. Many local maxima  
355 and minima are scattered around the parameter space, indicating that the  
356 surfaces behave very irregularly and exhibit highly nonlinear relationship with  
357 the input parameters. m sle, meters of sea level equivalent. 23

358 S4 Leave-one-out cross-validation results for the emulation performance. Each  
359 grey curve shows the comparison of zonal mean ice thickness transects from  
360 the model output and that from the emulator output for each parameter  
361 setting. Each boxplot shows the distribution of emulator output for each of  
362 the evenly spaced bins that span the range of true model output. In spite of  
363 the fact that our design points for parameter settings are quite sparse (100  
364 runs in 5-dimensional space) most of the curves are concentrated around 1:1  
365 line connecting the lower left and upper right corners of the plot, indicating  
366 that our emulator can reconstruct the original model output reasonably well  
367 across the input parameter settings.

24

368 S5 Leave-one-out cross-validation results for ice volume change projections across  
369 all 100 input parameter settings as the synthetic truth. The left panel shows  
370 95% prediction intervals for ice volume change projections across all 100 per-  
371 fect model experiments conducted for cross-validation. If the interval covers  
372 the 1:1 line connecting the lower left and upper right corners of the plot, the  
373 95% prediction interval includes the ice volume projection given by the syn-  
374 thetic truth. The right panel shows the coverage of those prediction intervals  
375 as a function of allowed range for the ice volume in 2005 AD relative to the  
376 observed ice volume. As going from left to right, the synthetic truths used in  
377 computing the coverages include more “unrealistic” ones in terms of modern  
378 ice volume. The numbers above the solid black line show how many synthetic  
379 truths fall into the given ice volume range. The plot shows that (i) the credible  
380 intervals achieve the nominal coverage level only for the “realistic” synthetic  
381 truths with modern ice volume within 10% of the observed ice volume, and  
382 (ii) the discrepancy term reduces overconfidence for the synthetic truths that  
383 are not within the 10% range.

25

384 S6 Summary of interactions between input parameters computed from leave-one-  
385 out cross-validation. Each panel shows the distribution of the rank correlation  
386 between two input parameters across all synthetic truths in our leave-one-out  
387 cross-validation. Five pairs of input parameters, (i) the flow factor and the  
388 snow PDD factor, (ii) the flow factor and the geothermal heat flux, (iii) the  
389 basal sliding factor and the ice PDD factor, (iv) the geothermal heat flux and  
390 the ice PDD, and (v) the ice PDD factor and the snow PDD factor are tend  
391 to be more negatively correlated comparing to the other pairs of parameters. 26

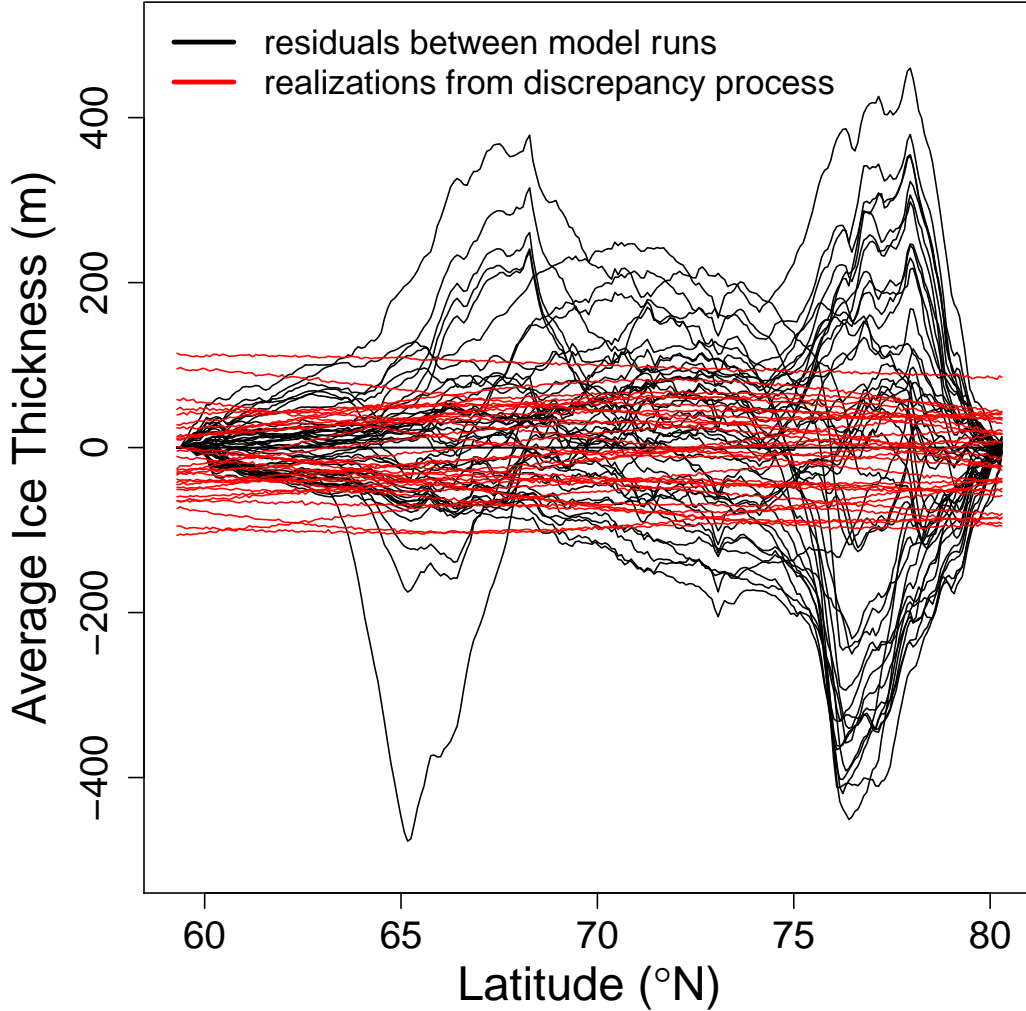


FIG. S1. Comparison between (i) residuals between the synthetic truth used in the main text (model run #67 in Applegate et al. 2012) and other model runs (black solid curves) and (ii) 30 different realizations from the model for the simulated discrepancy (red solid curves). The residuals are computed by subtracting the synthetic truth from each of the other model runs. For better display, we show only residual curves whose ranges are within  $(-500, 500)$ . It is easy to see that the black curves and red curves are generated from different processes, and therefore those two groups of curves can be separated by statistical inference (hence the discrepancy is identifiable). The magnitudes of the simulated discrepancy processes are well within the range covered by the model runs (hence the posterior density of input parameters does not show too large variation).

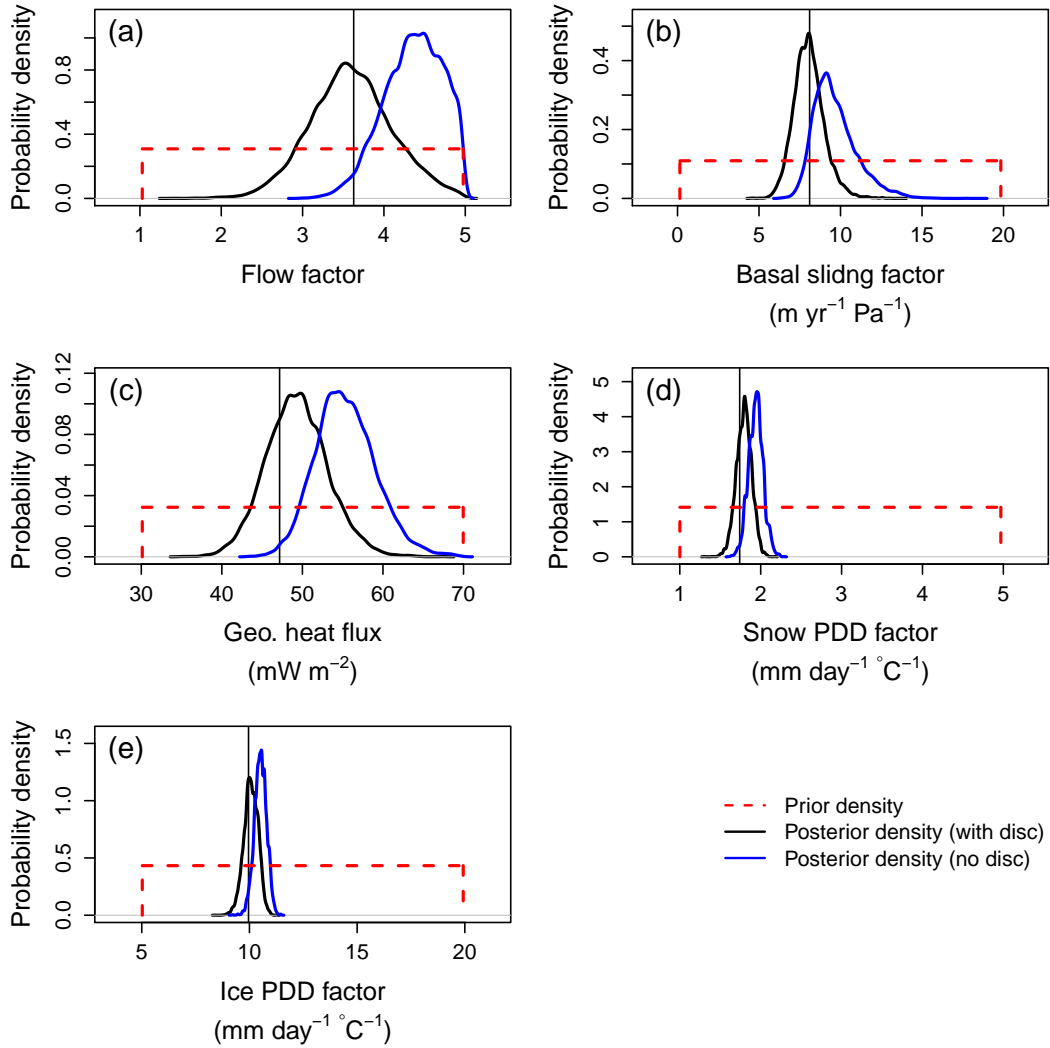


FIG. S2. Comparison between calibration results with and without the discrepancy term  $\mathbf{K}_d \boldsymbol{\nu}$  in the calibration model in (S1). In each panel, we tried to learn each of the parameters while fixing the other parameters at their assumed-true values. The prior densities are assumed to be uniform over a broad range (dashed red lines). While the posterior densities computed by including the discrepancy term in the model (solid black curves) pick up the true parameter values without notable biases, the posterior densities without the discrepancy term (solid blue curves) cannot recover the true values.

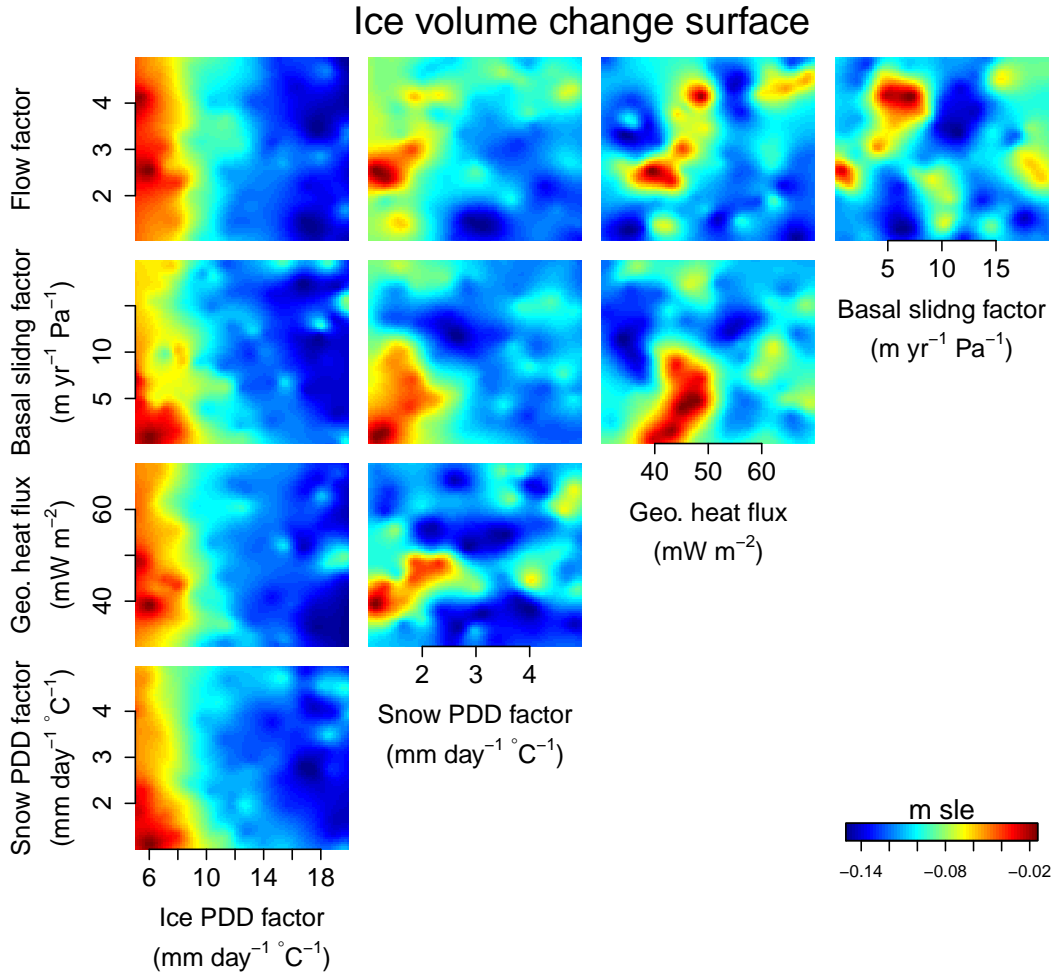


FIG. S3. Surfaces of ice volume change projections between 2005 and 2100 projected onto marginal spaces of all pairs of input parameters. Many local maxima and minima are scattered around the parameter space, indicating that the surfaces behave very irregularly and exhibit highly nonlinear relationship with the input parameters. m sle, meters of sea level equivalent.



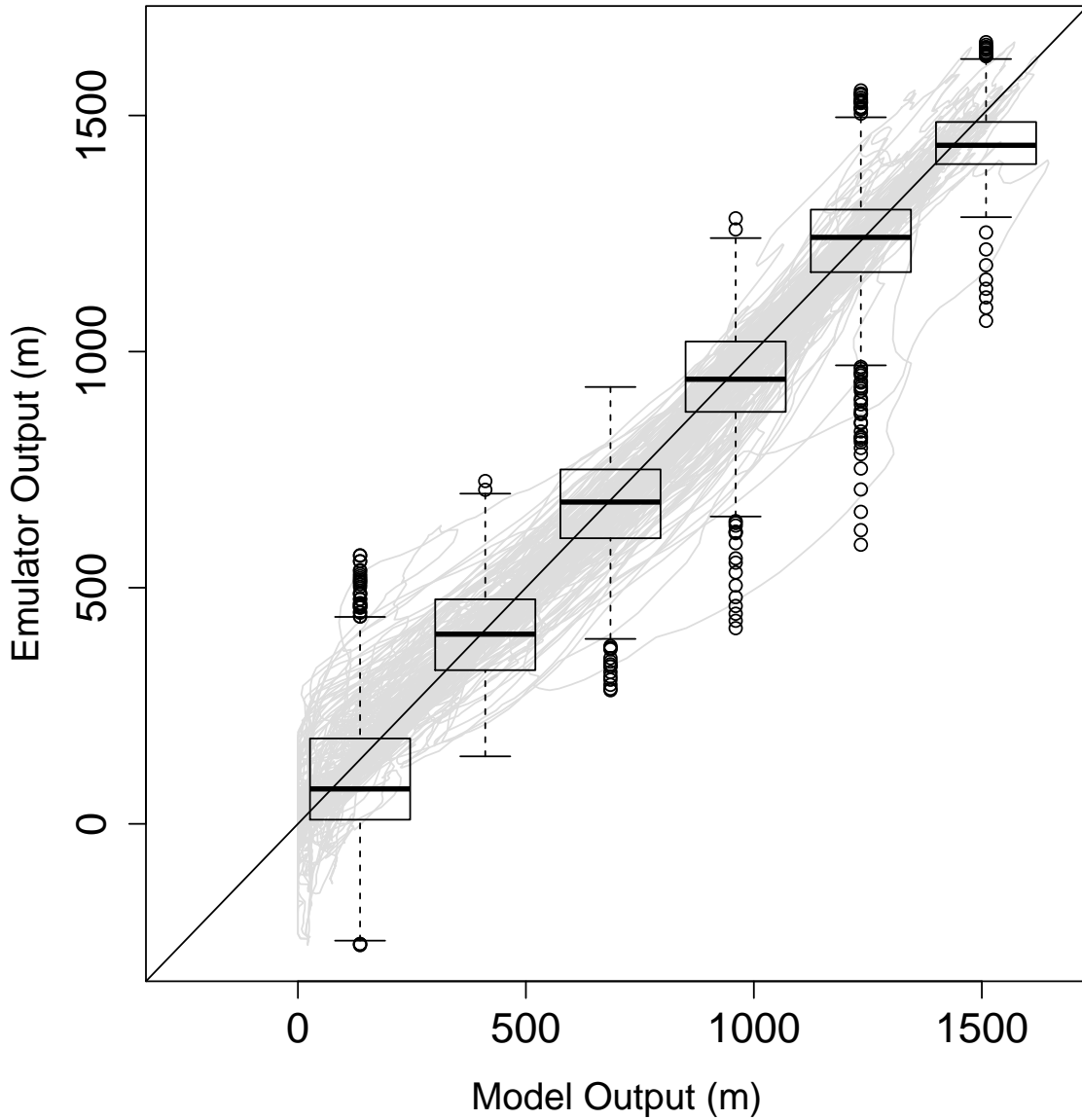


FIG. S4. Leave-one-out cross-validation results for the emulation performance. Each grey curve shows the comparison of zonal mean ice thickness transects from the model output and that from the emulator output for each parameter setting. Each boxplot shows the distribution of emulator output for each of the evenly spaced bins that span the range of true model output. In spite of the fact that our design points for parameter settings are quite sparse (100 runs in 5-dimensional space) most of the curves are concentrated around 1:1 line connecting the lower left and upper right corners of the plot, indicating that our emulator can reconstruct the original model output reasonably well across the input parameter settings.

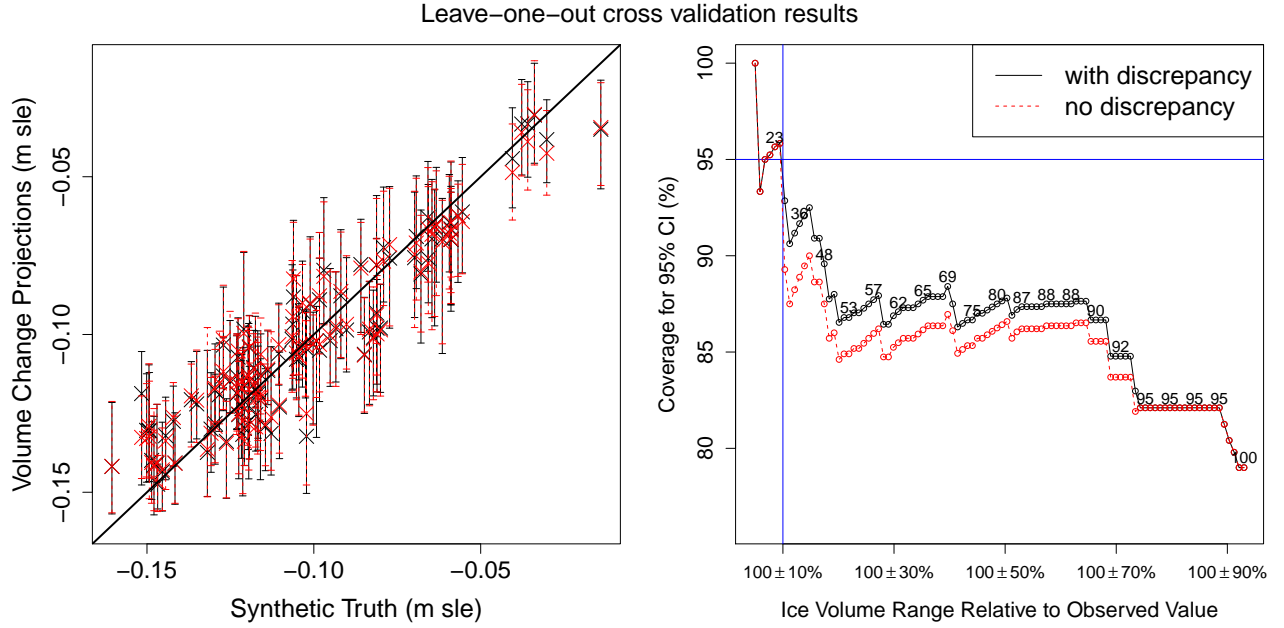


FIG. S5. Leave-one-out cross-validation results for ice volume change projections across all 100 input parameter settings as the synthetic truth. The left panel shows 95% prediction intervals for ice volume change projections across all 100 perfect model experiments conducted for cross-validation. If the interval covers the 1:1 line connecting the lower left and upper right corners of the plot, the 95% prediction interval includes the ice volume projection given by the synthetic truth. The right panel shows the coverage of those prediction intervals as a function of allowed range for the ice volume in 2005 AD relative to the observed ice volume. As going from left to right, the synthetic truths used in computing the coverages include more “unrealistic” ones in terms of modern ice volume. The numbers above the solid black line show how many synthetic truths fall into the given ice volume range. The plot shows that (i) the credible intervals achieve the nominal coverage level only for the “realistic” synthetic truths with modern ice volume within 10% of the observed ice volume, and (ii) the discrepancy term reduces overconfidence for the synthetic truths that are not within the 10% range.

Distributions of Rank Correlation between Input Parameters

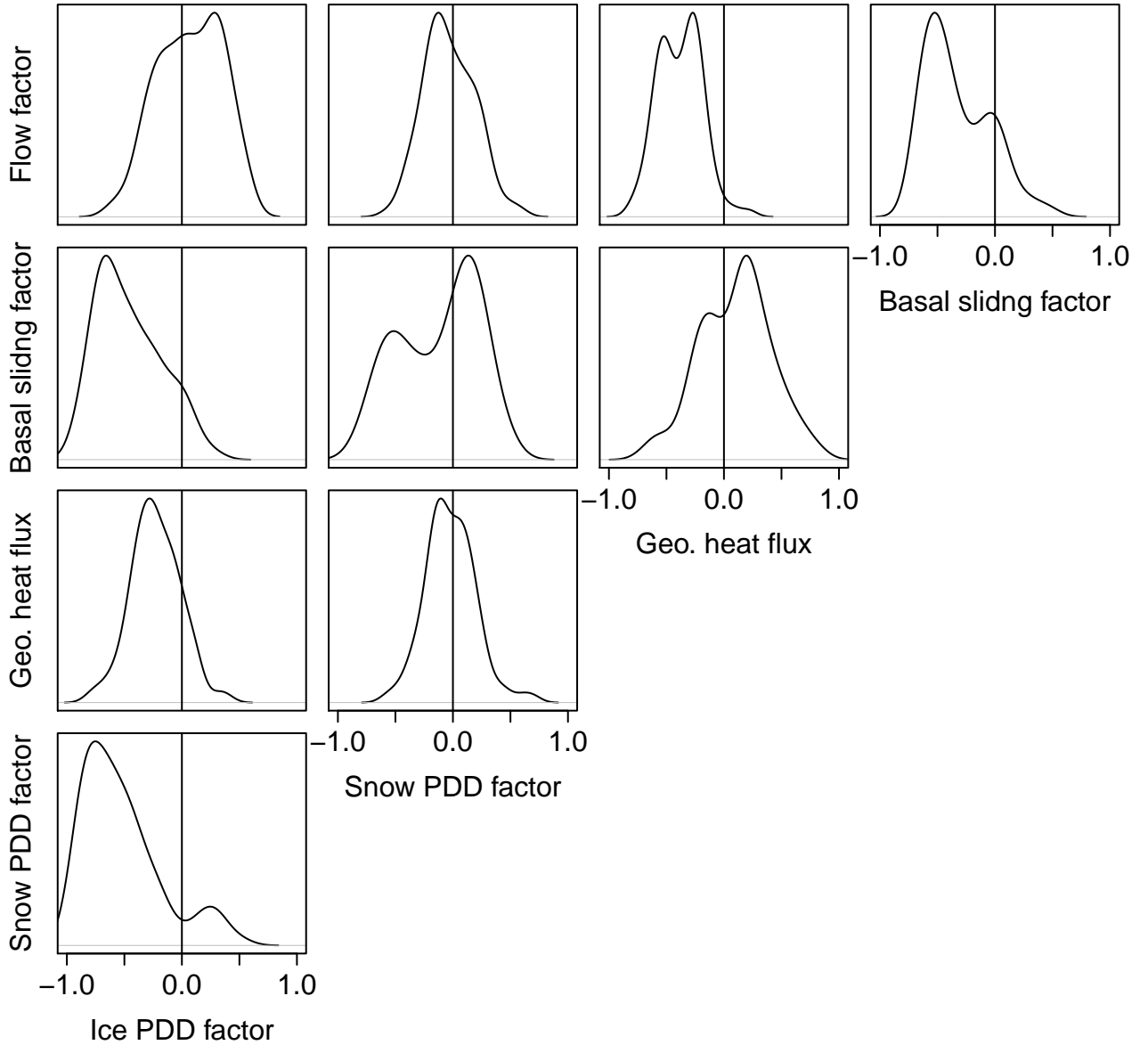


FIG. S6. Summary of interactions between input parameters computed from leave-one-out cross-validation. Each panel shows the distribution of the rank correlation between two input parameters across all synthetic truths in our leave-one-out cross-validation. Five pairs of input parameters, (i) the flow factor and the snow PDD factor, (ii) the flow factor and the geothermal heat flux, (iii) the basal sliding factor and the ice PDD factor, (iv) the geothermal heat flux and the ice PDD, and (v) the ice PDD factor and the snow PDD factor are tend to be more negatively correlated comparing to the other pairs of parameters.