

Supplement of Geosci. Model Dev., 11, 3131–3146, 2018
<https://doi.org/10.5194/gmd-11-3131-2018-supplement>
© Author(s) 2018. This work is distributed under
the Creative Commons Attribution 4.0 License.



Supplement of

Fast sensitivity analysis methods for computationally expensive models with multi-dimensional output

Edmund Ryan et al.

Correspondence to: Edmund Ryan (edmund.ryan@lancaster.ac.uk)

The copyright of individual parts of the supplement might differ from the CC BY 4.0 License.

1 **Supplementary Material**

2 **Section S1: Further Details of the Partial Least Squares approach**

3 PLS operates by projecting \mathbf{X} and \mathbf{Y} into new spaces, determined by maximising the covariance
4 between the projections of \mathbf{X} and \mathbf{Y} . The transformations of \mathbf{X} and \mathbf{X} are given by $\mathbf{X} = \mathbf{TP}^T + \mathbf{E}$
5 and $\mathbf{Y} = \mathbf{UQ}^T + \mathbf{F}$, where \mathbf{T} and \mathbf{U} are the projections of \mathbf{X} and \mathbf{Y} respectively, \mathbf{P} and \mathbf{Q} are
6 orthogonal matrices, and \mathbf{E} and \mathbf{F} are independent and identically distributed error terms. It is
7 straightforward to show that $\mathbf{Y} = \mathbf{XB} + \mathbf{G}$, where $\mathbf{B} = \mathbf{X}^{-1}\mathbf{TQ}^T$ and $\mathbf{G} = (\mathbf{U} - \mathbf{T})\mathbf{Q}^T + \mathbf{F}$. The
8 \mathbf{B} matrix stores the PLS-regression coefficients, which can be interpreted as sensitivity indices
9 (Chang et al., 2015; Sobie, 2009).

10 **Section S2: Setting up the global sensitivity analysis experiment**

11 Prior to running the scripts accompanying this paper, the following packages are required:

```
12 install.packages("lhs"); install.packages("emulator"); install.packages("mvtnorm");  
13 install.packages("mgcv"); install.packages("sensitivity"); install.packages("DiceKriging");  
14 install.packages("DiceOptim")
```

15 (For first time users of the open source programming language R, go to [https://www.r-](https://www.r-project.org/)
16 [project.org/](https://www.r-project.org/) to download it for free). We decided on the model inputs/parameters, the ranges of
17 the inputs, and the outputs. For applications to other problems, note that the outputs could be a
18 spatial map, a time-series or a combination of both. We created the design matrix by running the
19 following as a new script file, specifying the number of input factor and their ranges:

```
20 #Things to specify by the user:  
21 setwd("C:/Users/...") #Location of folder where files are stored.  
22 Np = ?? #No. of input factors
```

```

1  mink = c(p1min,p2min,p3min,...) #Min values of the inputs/parameters.
2  maxk = c(p1max,p2max,p3max,...) #Max values of the inputs/parameters.
3  library(lhs); library(emulator)
4  inputs_norm = maximinLHS(Np*10,Np)
5  write.table(inputs_norm,"InputsNorm_TrainingData.csv",row.names=F,col.names=F,sep=",")
6  inputs = matrix(-9999,nrow=k*10,ncol=k)
7  for (i in 1:k){inputs[,i] = (Inputs_norm[,i]*(maxk[i]-mink[i])) + mink[i]}c
8  write.table(inputs,"Inputs_TrainingData.csv", row.names=F, col.names=F, sep=",")

```

9 For input ranges that were on the log-scale (e.g. the min/max of input/parameter p is $0.01 * pCtrl$
10 and $100 * pCtrl$, where $pCtrl$ is the control run value of the input/parameter p), then we first
11 transformed to a linear scale before running the script. We ran the chemistry model for each of
12 the rows of the *Inputs_TrainingData.csv* file and stored the outputs in a csv file
13 *Outputs_TrainingData.csv*. The outputs in this csv file consisted of N_x rows and N_y columns,
14 where $N_x = N_p * 10$ is the number of runs of the model, and N_y is the length of the row vector
15 storing the output for a given input.

16 **Section S3: Calculating first order and total sensitivity indices using the Sobol method** 17 **without dimension reduction**

18 We run the following R script file:

```

19 X=read.csv('InputsNorm_TrainingData.csv', header=FALSE)
20 yALL =read.csv('Outputs_TrainingData.csv', header=FALSE)
21 SI = matrix(-9999,nrow=dim(inputs)[2],ncol= dim(outputs)[2])
22 SI.total = matrix(-9999,nrow=dim(inputs)[2],ncol= dim(outputs)[2])

```

```

1  for (j in 1: dim(outputs)[2]){y = as.matrix(yALL[,j],rownames.force=NA)
2  m = km(~ ., design = X, response = y, covtype = "matern3_2")
3  AB = randomLHS(N,16); A = AB[,1: dim(X)[2]]; B = AB[, (dim(X)[2]+1):(dim(X)[2]*2)]
4  yA = predict.km(m, A, "UK", se.compute = FALSE, checkNames = FALSE)$mean
5  yB = predict.km(m, B, "UK", se.compute = FALSE, checkNames = FALSE)$mean
6  yA_sum = mean(yA); yB_sum = mean(yB); yAyA_sum = mean(yA^2); yByB_sum = mean(yB^2)
7  for (i in 1: dim(X)[2]){print(c(j,i)); Ci=B; Ci[,i]=A[,i]; Di=A; Di[,i]=B[,i]
8  yCi = predict.km(m, Ci, "UK", se.compute = FALSE, checkNames = FALSE)$mean
9  yDi = predict.km(m, Di, "UK", se.compute = FALSE, checkNames = FALSE)$mean
10 yCi_sum = mean(yCi); yDi_sum = mean(yDi); yCiyCi_sum = mean(yCi^2);
11 yDiyDi_sum = mean(yDi^2); yAyCi_sum = mean(yA*yCi);
12 yAyDi_sum = mean(yA*yDi); yByCi_sum = mean(yB*yCi);
13 yByDi_sum = mean(yB*yDi); SI_temp = rep(-9999, 8)
14 SI_temp[1] = (yAyCi_sum - (yA_sum*yB_sum))/(yAyA_sum - (yA_sum*yB_sum))
15 SI_temp[2] = (yByDi_sum - (yA_sum*yB_sum))/(yByB_sum - (yA_sum*yB_sum))
16 SI_temp[3] = (yAyCi_sum - (yA_sum*yB_sum))/(yByB_sum - (yA_sum*yB_sum))
17 SI_temp[4] = (yAyCi_sum - (yCi_sum*yDi_sum))/(yCiyCi_sum - (yCi_sum*yDi_sum))
18 SI_temp[5] = (yAyCi_sum - (yCi_sum*yDi_sum))/(yDiyDi_sum - (yCi_sum*yDi_sum))
19 SI_temp[6] = (yByDi_sum - (yA_sum*yB_sum))/(yAyA_sum - (yA_sum*yB_sum))
20 SI_temp[7] = (yByDi_sum - (yCi_sum*yDi_sum))/(yCiyCi_sum - (yCi_sum*yDi_sum))
21 SI_temp[8] = (yByDi_sum - (yCi_sum*yDi_sum))/(yDiyDi_sum - (yCi_sum*yDi_sum))
22 SI[i,j] = mean(SI_temp)}}
23 write.table(SI*100, "SIs_Sobol_EmulatorOnly.csv", row.names=F, col.names=F, sep=",")

```

1 To compute the total indices, the eight SI_temp lines can be replaced with the following:

```
2 SI_temp[1] = (yAyDi_sum - (yA_sum*yB_sum))/(yAyA_sum - (yA_sum*yB_sum))
3 SI_temp[2] = (yByCi_sum - (yA_sum*yB_sum))/(yByB_sum - (yA_sum*yB_sum))
4 SI_temp[3] = (yAyDi_sum - (yCi_sum*yDi_sum))/(yAyA_sum - (yCi_sum*yDi_sum))
5 SI_temp[4] = (yByCi_sum - (yCi_sum*yDi_sum))/(yByB_sum - (yCi_sum*yDi_sum))
```

6 **Section S4: Calculating first order and total sensitivity indices using the Sobol method,**
7 **using principal component analysis to reduce the dimensionality of the output**

8 We run the following R script file to compute the sensitivity indices for the emulator-PCA
9 hybrid approach. The R script below can be modified to ensure that the number of principal
10 components included accounts for a specified proportion of the variance (99% used here).

```
11 X=read.csv('InputsNorm_TrainingData.csv', header=FALSE)
12 yALL =read.csv('Outputs_TrainingData.csv', header=FALSE)
13 S = var(Y); S.eig = eigen(S)
14 eig.value.cumul = cumsum((S.eig$values/sum(S.eig$values))*100)
15 yALL.dim.index = c(1:dim(yALL)[2]); Npca=min(yALL.dim.index[eig.value.cumul>99])
16 ptm <- proc.time()
17 AB <- randomLHS(N,(dim(X)[2]*2))
18 A <- AB[,1:dim(X)[2]]; B <- AB[(dim(X)[2]+1):(dim(X)[2]*2)]
19 yA_PCA=matrix(-9999,N,Npixels); yB_PCA=matrix(-9999,N,Npixels)
20 for (j in 1:Npc){
21 PC <- matrix(rep(S.eig$vectors[,j],dim(X)[1]),nrow=dim(X)[1],ncol=Npixels,byrow=TRUE)
22 y_PCj <- as.matrix(rowSums(PC*Y),rownames.force=NA)
```

```

1  assign(paste("m",j, sep=""),km(~., design = X, response = y_PCj, covtype = "matern5_2"))
2  yA_PCA[,j] <- predict.km(eval(as.symbol(paste("m",j,sep=""))), A, "UK", se.compute =
3  FALSE, checkNames = FALSE)$mean
4  yB_PCA[,j] <- predict.km(eval(as.symbol(paste("m",j,sep=""))), B, "UK", se.compute =
5  FALSE, checkNames = FALSE)$mean
6  }
7  for (j in (Npc+1):Npixels){
8  yA_PCA[,j]=matrix(0,nrow=N,ncol=1); yB_PCA[,j]=matrix(0,nrow=N,ncol=1)
9  }
10 PCs.inv <- solve(S.eig$vector); yA <- yA_PCA%%PCs.inv; yB <- yB_PCA%%PCs.inv
11 yCi <- matrix(-9999,nrow=N,ncol=Npixels*dim(X)[2])
12 yDi <- matrix(-9999,nrow=N,ncol=Npixels*dim(X)[2])
13 for (i in 1:dim(X)[2]){
14  print(c(i)); Ci=B; Ci[,i]=A[,i]; Di=A; Di[,i]=B[,i]
15  yCj_PCA <- matrix(-9999,nrow=N,ncol=Npixels)
16  yDj_PCA <- matrix(-9999,nrow=N,ncol=Npixels)
17  for (j in 1:Npc){
18  yCj_PCA[,j] <- predict.km(eval(as.symbol(paste("m",j,sep=""))), Ci, "UK", se.compute =
19  FALSE, checkNames = FALSE)$mean
20  yDj_PCA[,j] <- predict.km(eval(as.symbol(paste("m",j,sep=""))), Di, "UK", se.compute =
21  FALSE, checkNames = FALSE)$mean
22  }
23  for (j in (Npc+1):Npixels){

```

```

1  yCj_PCA[,j] <- matrix(0,nrow=N,ncol=1); yDj_PCA[,j] <- matrix(0,nrow=N,ncol=1)
2  }
3  i2 <- i*Npixels; i1 <- i2 - (Npixels-1)
4  yCi[,i1:i2] <- yCj_PCA%%PCs.inv; yDi[,i1:i2] <- yDj_PCA%%PCs.inv
5  }
6  SI <- matrix(-9999,nrow=8,ncol=Npixels)
7  for (i in 1:Npixels){
8    yA_sum <- mean(yA[,i]); yB_sum <- mean(yB[,i]);
9    yAyA_sum <- mean(yA[,i]^2); yByB_sum <- mean(yB[,i]^2)
10   SI_temp <- rep(-9999,8)
11   for (j in 1:8){
12     z <- ((j-1)*Npixels) + 1; yCi_sum <- mean(yCi[,z]); yDi_sum <- mean(yDi[,z])
13     yCiyCi_sum <- mean(yCi[,z]^2); yDiyDi_sum <- mean(yDi[,z]^2)
14     yAyCi_sum <- mean(yA[,i]*yCi[,z]); yAyDi_sum <- mean(yA[,i]*yDi[,z])
15     yByCi_sum <- mean(yB[,i]*yCi[,z]); yByDi_sum <- mean(yB[,i]*yDi[,z])
16     SI_temp[1] <- (yAyCi_sum - (yA_sum*yB_sum))/(yAyA_sum - (yA_sum*yB_sum))
17     SI_temp[2] <- (yByDi_sum - (yA_sum*yB_sum))/(yByB_sum - (yA_sum*yB_sum))
18     SI_temp[3] <- (yAyCi_sum - (yA_sum*yB_sum))/(yByB_sum - (yA_sum*yB_sum))
19     SI_temp[4] <- (yAyCi_sum - (yCi_sum*yDi_sum))/(yCiyCi_sum - (yCi_sum*yDi_sum))
20     SI_temp[5] <- (yAyCi_sum - (yCi_sum*yDi_sum))/(yDiyDi_sum - (yCi_sum*yDi_sum))
21     SI_temp[6] <- (yByDi_sum - (yA_sum*yB_sum))/(yAyA_sum - (yA_sum*yB_sum))
22     SI_temp[7] <- (yByDi_sum - (yCi_sum*yDi_sum))/(yCiyCi_sum - (yCi_sum*yDi_sum))
23     SI_temp[8] <- (yByDi_sum - (yCi_sum*yDi_sum))/(yDiyDi_sum - (yCi_sum*yDi_sum))

```

```

1  SI[j,i] <- mean(SI_temp)
2  } }
3  proc.time() - ptm #Stop the clock
4  write.csv(data.frame(SI*100), file =
5  paste("SIs_Sobol_EmulatorPCA_",modelN,"_N",N,".csv",sep=""), row.names=FALSE)

```

6 Section S5: Calculating first order and total sensitivity indices using the extended FAST 7 (eFAST) method

8 We ran the following as a new R script file:

```

9  library(sensitivity); library(DiceKriging); library(DiceOptim)
10 X=read.csv('InputsNorm_TrainingData.csv', header=FALSE)
11 yALL =read.csv('Outputs_TrainingData.csv', header=FALSE)
12 SI = matrix(-9999,nrow=dim(X)[2],ncol= dim(yALL)[2])
13 SI.total = matrix(-9999,nrow=dim(X)[2],ncol= dim(yALL)[2])
14 for (j in 1: dim(yALL)[2]){y = as.matrix(yALL[,j],rownames.force=NA)
15 m = km(~ ., design = X, response = y, covtype = "matern3_2")
16 kriging.mean = function(Xnew, m){predict.km(m, Xnew, "UK", se.compute = FALSE,
17 checkNames = FALSE)$mean}
18 temp = fast99(model = kriging.mean, factors = dim(X)[2], n = 1000, q = "qunif", q.arg =
19 list(min = 0, max = 1), m = m)
20 SI[,j] = as.matrix(temp$D1/temp$V)
21 SI.total[,j] = as.matrix((temp$V- temp$Dt)/temp$V)}
22 write.table(SI*100," SIs_eFAST_EmulatorOnly.csv", row.names=F, col.names=F, sep=",")

```



```
1 write.table(SI.total*100, "TotalSIs_eFAST_EmulatorOnly.csv", row.names=F, col.names=F,  
2 sep=";",")
```

3 **Section S6: Basic emulator diagnostic checks**

4 We run the following R script file:

```
5 library(sensitivity); library(DiceKriging); library(DiceOptim); library(lhs);  
6 X=read.csv('InputsNorm_TrainingData.csv', header=FALSE)  
7 inputs_val_norm = maximinLHS(30,dim(X)[2])  
8 write.table(inputs_val_norm,"InputsNorm_TrainingData.csv", row.names=F, col.names=F,  
9 sep=";",")  
10 inputs_val = matrix(-9999,nrow=30,ncol=dim(X)[2])  
11 for (i in 1:dim(X)[2]){inputs[,i] = (inputs_val_norm[,i]*(max(X[,i])-min(X[,i])) + min(X[,i]  
12 write.table(inputs_val,"Inputs_ValidationRuns.csv", row.names=F, col.names=F, sep=";",")
```

13 Each row of the *Inputs_ValidationRuns.csv* file provides the input values for a given model
14 simulation, and we write the output to a corresponding row in the output file
15 *Outputs_ValidationRuns.csv*. The diagnostic check compares the outputs of the simulator at the
16 validation inputs with those predicted by the emulator. We run the following R script file to
17 carry out this diagnostic check:

```
18 X=read.csv('InputsNorm_TrainingData.csv', header=FALSE)  
19 yALL =read.csv('Outputs_TrainingData.csv', header=FALSE)  
20 X_val=read.csv('InputsNorm_ValidationRuns.csv', header=FALSE)  
21 yALL_val =read.csv('Outputs_ValidationRuns.csv', header=FALSE)  
22 for (j in 1: dim(yALL)[2]){y = as.matrix(yALL[,j],rownames.force=NA)
```

```

1  m = km(~ ., design = X, response = y, covtype = "matern5_2")
2  yALL_val_emul[,j] = as.matrix(predict.km(m, X_val, "UK", se.compute=FALSE,
3  checkNames=FALSE)$mean)}
4  plot(c(yALL_val),c(yALL_val_emul),main="Emulator versus Simulator predictions at validation
5  inputs",xlab="simulator", ylab="emulator").
6  print(cor(c(x),c(y))^2)*100 #This prints the value of the coefficient of determination R2.

```

7 **Section S7: Calculating first and higher order sensitivity indices using the Generalized**
8 **Additive Model (GAM) method.**

9 We ran the following as a new R script file:

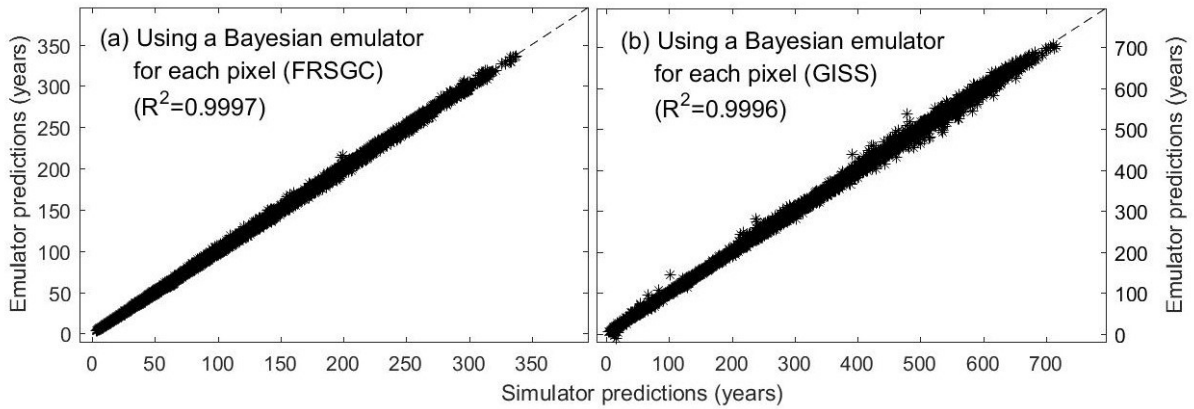
```

10 library(mgcv);
11 X = read.csv('InputsNorm_TrainingData.csv', header=FALSE)
12 yALL = read.csv('OutputsNorm_TrainingData.csv', header=FALSE)
13 SI = matrix(-9999,nrow=dim(inputs)[2],ncol= dim(outputs)[2])
14 SI.total = matrix(-9999,nrow=dim(inputs)[2],ncol= dim(outputs)[2])
15 for (j in 1:dim(outputs)[2]){y = as.matrix(yALL[,j],rownames.force=NA)
16 SI = matrix(-9999,nrow=dim(inputs)[2],ncol=dim(outputs)[2])
17 for (j in 1:dim(outputs)[2]){y=yALL[,j]; vary=var(Y); v=rep(-9999,8);
18 for (i in 1:dim(inputs)[2]){gam.model = gam(Y ~ te(X[,i])); v[i]=var(gam.model$fitted)}-
19 SI[,j]=(v/varY)*100}
20 write.table(SI, "SIs_GAM.csv", row.names=F, col.names=F, sep=", ")

```

21 Note that the second order SIs of the *i*th and *k*th inputs can be computed by replacing *gam(Y ~*
22 *te(X[,i]))* with *gam.model = gam(Y ~ ti(X[,i], X[,k]))*.

1



2

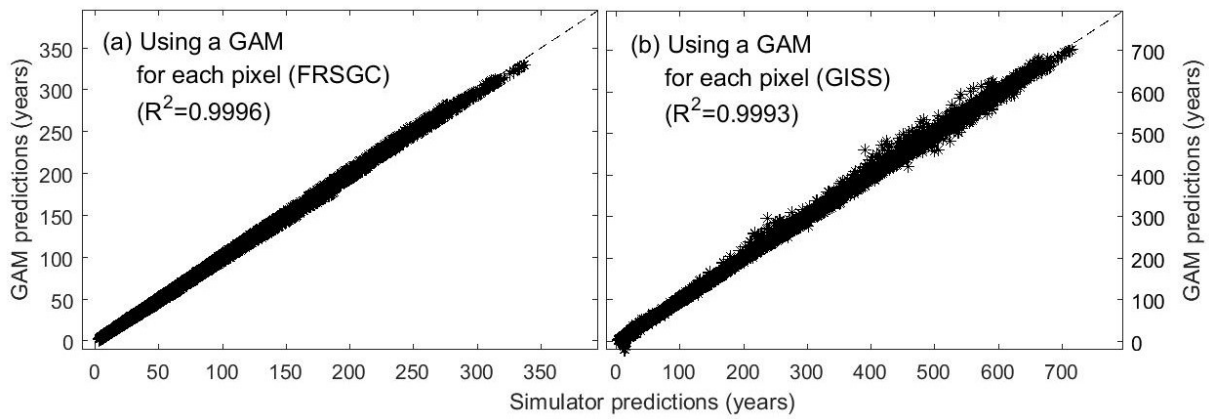
3

Figure S1. Annual column mean CH₄ lifetime calculated by the FRSGC and GISS chemistry models (x-axis) versus predicted by the Gaussian Process emulator (y-axis).

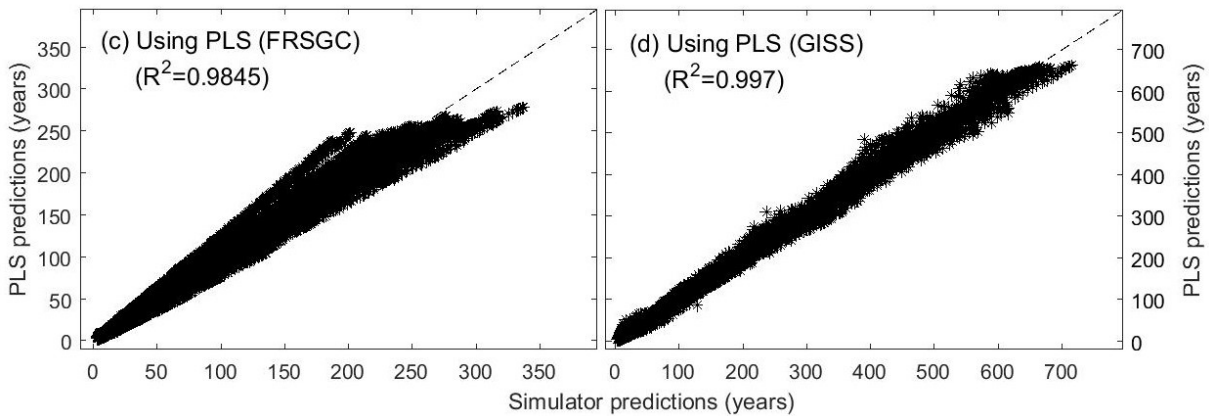
4

5

6



7



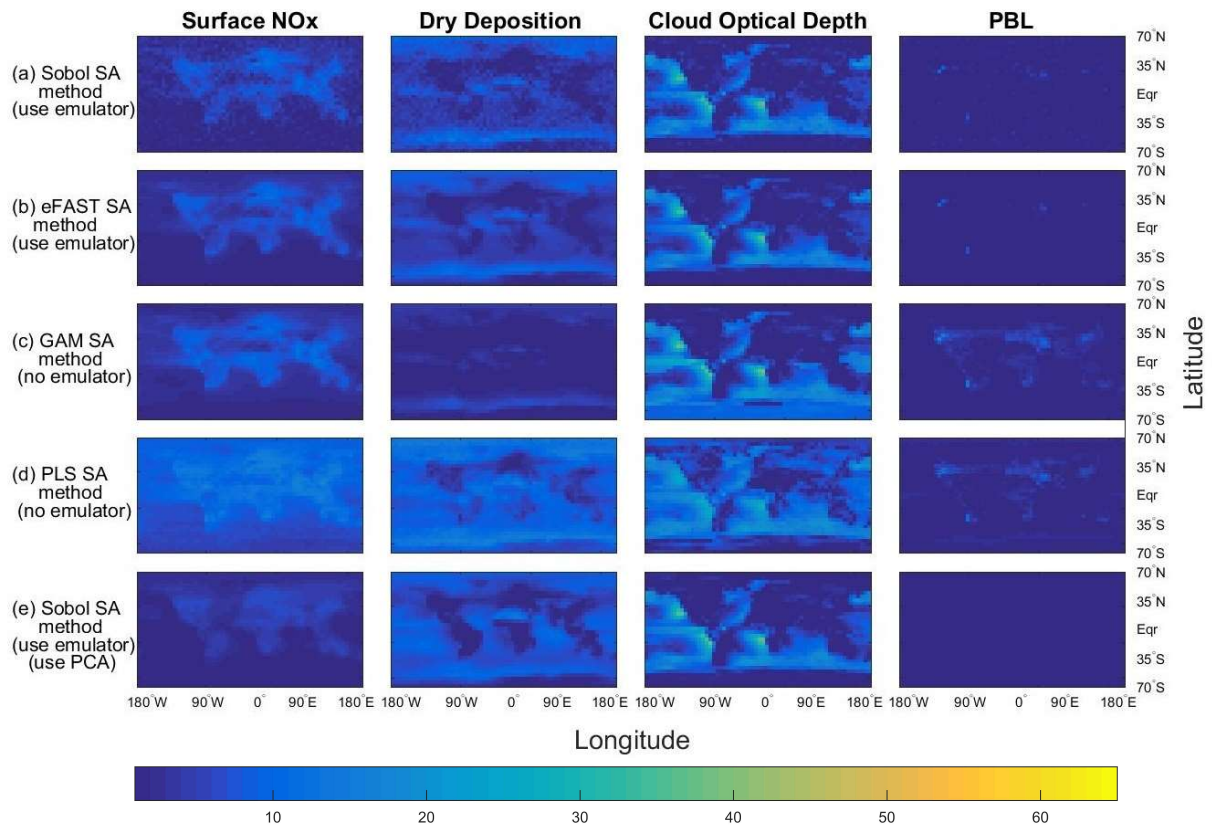
8

9

10

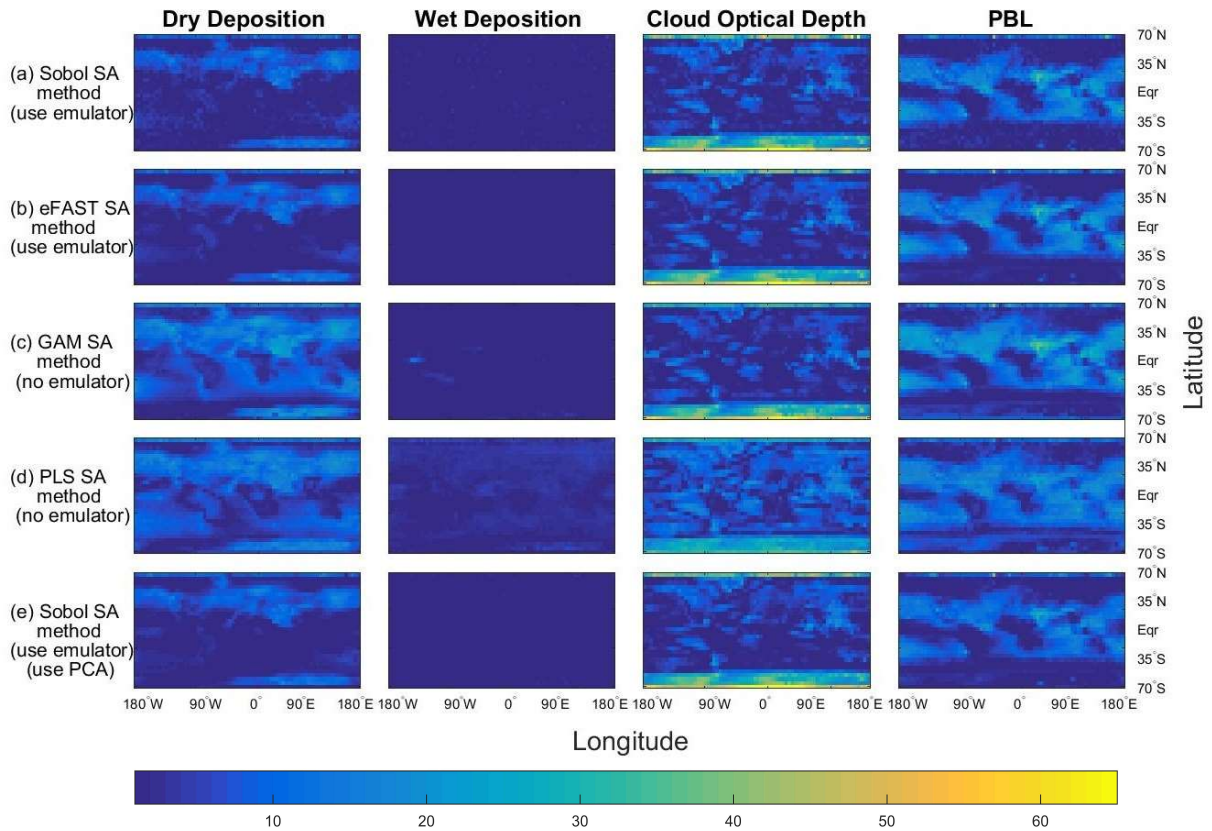
11

Figure S2. Annual column mean CH₄ lifetime calculated by the FRSGC and GISS chemistry models (x-axis) versus predictions by a separate generalized additive model (panels a and b) or partial least squares model (panels c and d) for each pixel (y-axis).



1
2
3
4
5

Figure S3. Sensitivity indices for the four minor inputs for the FRSGC chemistry transport model showing the inputs not already included in Figure 3.



1

2

3 **Figure S4.** Sensitivity indices for the four minor inputs for the GISS chemistry transport model
 4 showing the inputs not already included in Figure 4.