

# Introduction to emulators

## - the what, the when, the why

Dr Lindsay Lee

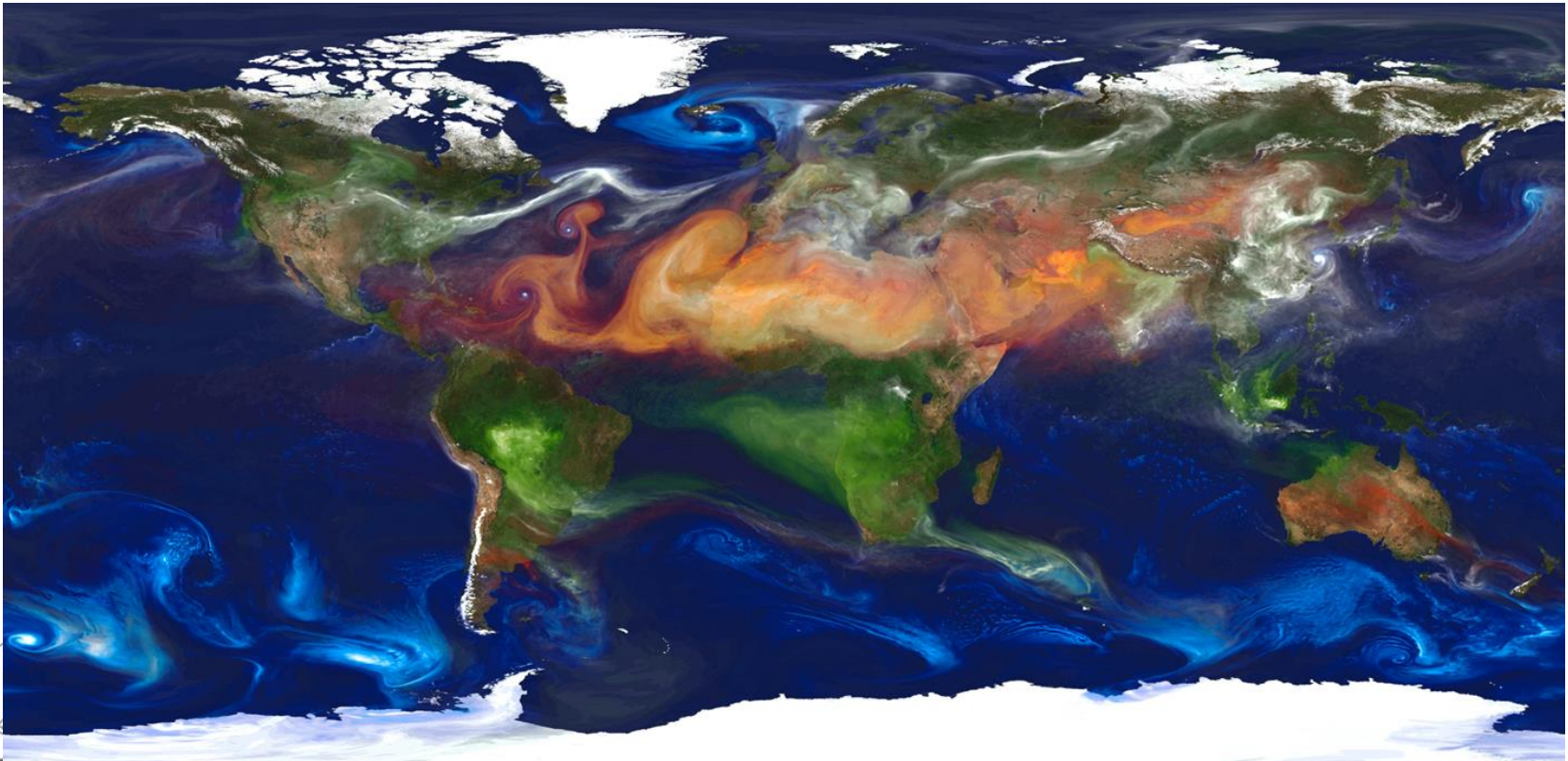


# What is a simulator?



UNIVERSITY OF LEEDS

A simulator is a computer code used to represent some real world process



## Understand

Which aerosol  
are most  
effective at  
cooling the  
climate?

## Predict

Could aerosol be  
used to cool the  
climate for  
geoengineering?



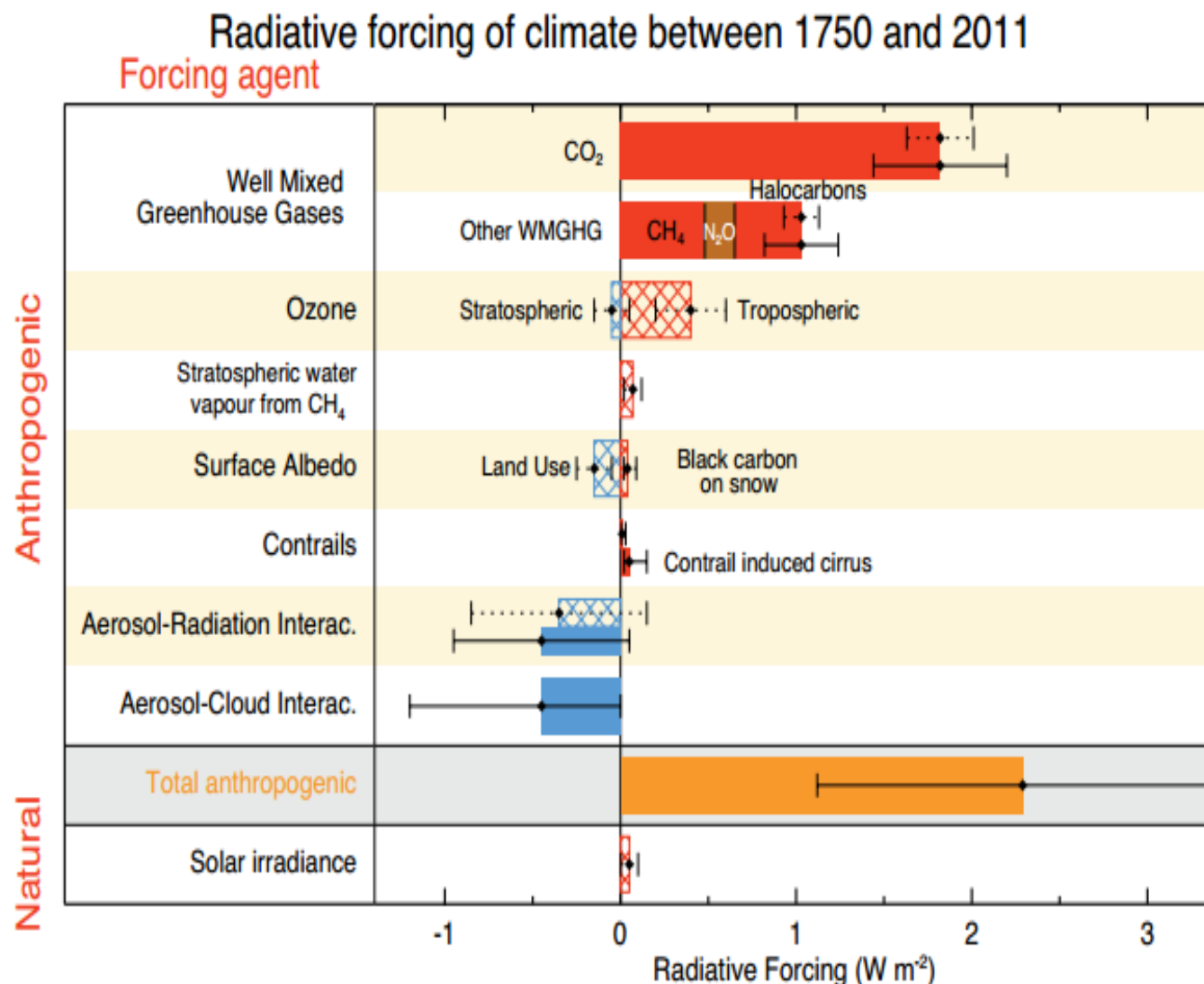
- $y = f(x)$  represents the simulation process
- $y$  is the simulator output
- $x$  is the simulator input
- $f$  is the simulator



$$Y = f(X)$$

Computer codes are  
imperfect representations of  
the real world.

How do these imperfections affect our  
understanding and predictions?  
What is the effect on  $Y$ ?



**Figure 8.15** | Bar chart for RF (hatched) and ERF (solid) for the period 1750–2011, where the total ERF is derived from Figure 8.16. Uncertainties (5 to 95% confidence range) are given for RF (dotted lines) and ERF (solid lines).



- $X$  are 'model inputs/parameters'
- They may be spatial fields/time series
- They may be single values used globally or regionally

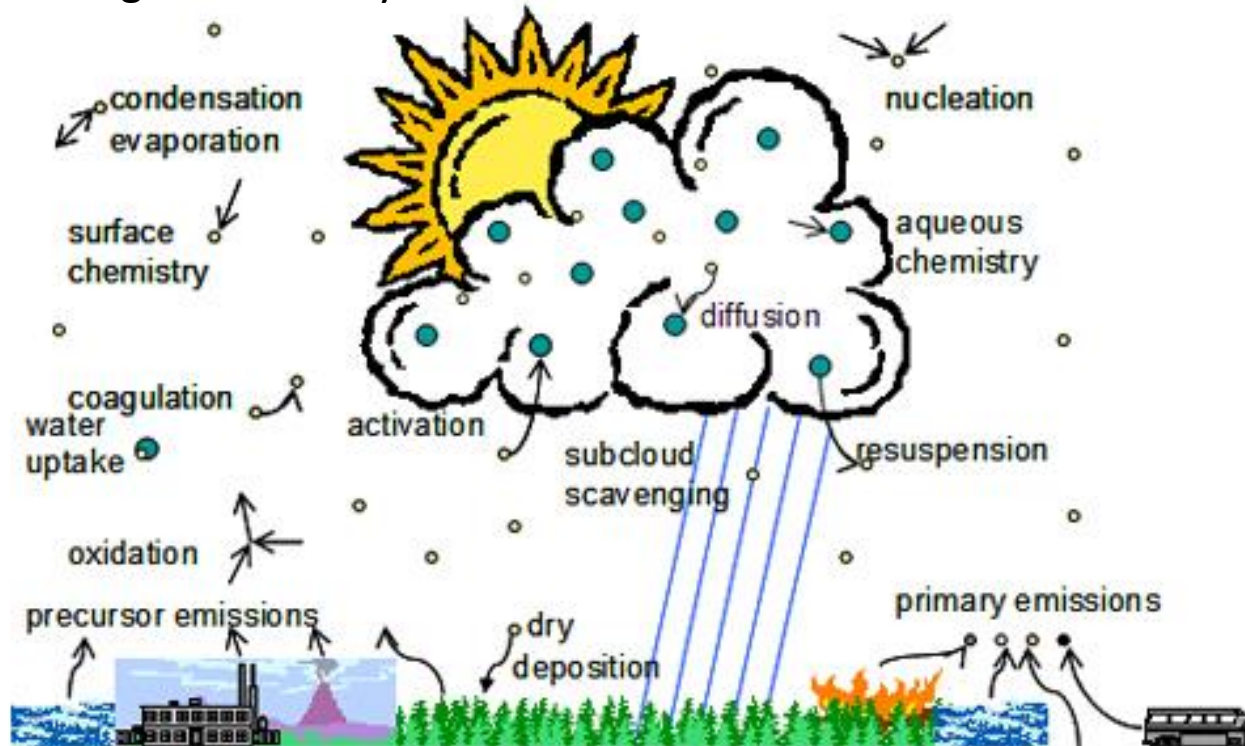


We will represent the uncertainty in  $X$  by a single value to perturb the whole field, or series (maybe regionally, but mostly globally)





Diagram courtesy of PNNL

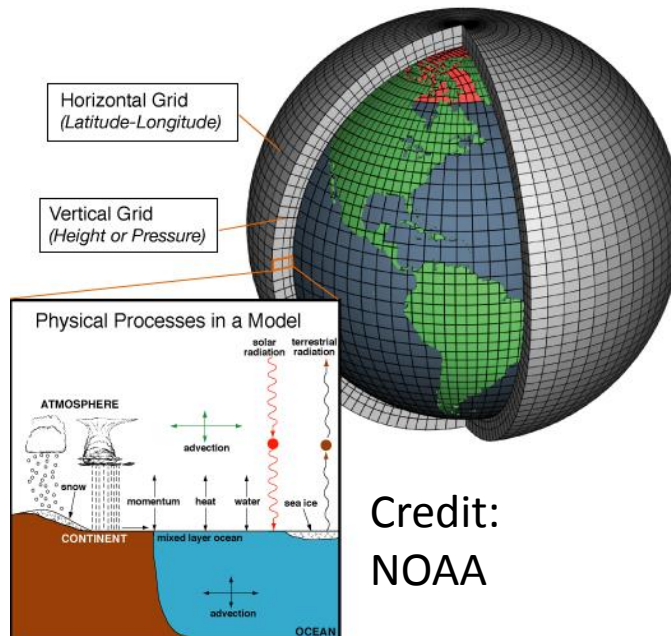


- $X$  is uncertain  $\rightarrow Y$  is uncertain
- $x$  can be given uncertainty distribution  $G$
- Marginally  $x_k$  has distribution  $G_k$

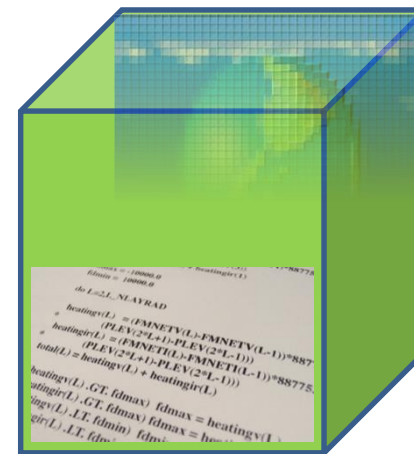
Obtain  $y_i = f(x_i)$  for given  $x_i$   
from  $G$



- Simulators contain lots of equations (hundreds of lines of code)
- Simulator resolution increases with computer power



Credit:  
NOAA



$f$  is often too complex to simulate  
many  $x_i$  from  $G$  to characterise the  
effect of uncertainty on  $y$



- Replace  $f$  with  $\hat{f}$
- $\hat{f}$  is the emulator
- $\hat{f}$  is much quicker to run but accurately represents  $f$



- Replace  $y_k = f(x)$  with  $y_k = \hat{f}(x)$

Define an output (or selection of) to be emulated

- doesn't replace the whole simulator for all research



- From the simulator  $y$  can be any/all model output/diagnostics
- For emulation  $y$  must be more restricted
- $y$  is often a scalar representing a single model output, at a particular time and location
  - Eg.
    - global mean temperature for 2000
    - total number for particulate matter aerosol  $> 2.5\mu\text{m}$  in a model grid box
- $y$  could be a spatial field, or a time series, or EOF/PCAs, or a selection of model outputs



The complexity of the emulator is dependent on the complexity of the model output to be emulated





- We will carry on assuming  $y$  is scalar and  $\mathbf{x}$  is a collection of scalar values
- $\mathbf{x} = x_1, x_2, x_3, \dots, x_p$  for  $p$  uncertain parameters
- Use a selection of  $\mathbf{x}_i, \mathbf{y}_i$  from  $f$  to 'build'  $\hat{f}$



# When is an emulator useful?



UNIVERSITY OF LEEDS

## uncertainty analysis

Quantify the impact of uncertainty on model output

## sensitivity analysis

Understand model response to uncertain inputs

And the model is too complex to do the sampling required



An emulator that gives estimates of its own uncertainty

We tend to use the Gaussian process



- Each point in space has a Gaussian distribution
- Each collection of points has a multivariate Gaussian distribution
- The whole collection is a Gaussian process
- In its basic form requires 'smooth output'

Does not require the output to be  
Gaussian



Set up the Gaussian process prior function

$$f(\mathbf{x}) | \beta, \sigma \sim GP(h(\mathbf{x})^T \beta, \sigma^2 c(\mathbf{x}, \mathbf{x}'))$$

$h(\mathbf{x})^T \beta$  is the mean function

$c(\mathbf{x}, \mathbf{x}')$  is the covariance function

$\beta, \sigma$  are hyperparameters



- Use the mean function to specify any functional form known to exist
- In the absence of prior information we tend to specify it as constant or a linear regression formula

$$h(x) = 1$$

$$h(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$



- When there are very few simulator runs to give information the function is weighted towards the prior mean
- Outside the bounds of the simulator runs the emulator tends towards the prior mean
- With more (well-chosen) simulator runs the hyperparameters  $\beta$  are better estimated



- Can be plugged in if a particular regression model is required
- Usually use maximum likelihood to estimate them and ‘plug-in’
- Could use a Bayesian approach but time-consuming





- The covariance function should reflect  $c(\mathbf{x}, \mathbf{x}) = 1$  and that  $c(\mathbf{x}, \mathbf{x}')$  decreases as  $|\mathbf{x} - \mathbf{x}'|$  increases
- Includes further hyperparameters  $\delta$  that specify the speed of covariance decrease with increasing  $|\mathbf{x} - \mathbf{x}'|$

$$\text{cov}\{f(\mathbf{x}), f(\mathbf{x}') | \sigma\} = \sigma^2 c(\mathbf{x}, \mathbf{x}')$$



**Gaussian:**

$$c(x, x') | \delta = \exp\left\{-\frac{(x - x')^2}{2\delta^2}\right\}$$

**Matern 5/2:**

$$c(x, x') | \delta = \left(1 + \frac{\sqrt{5}|x-x'|}{\delta} + \frac{5(x-x')^2}{3\delta^2}\right) \exp\left(-\frac{\sqrt{5}|x-x'|}{\delta}\right)$$

- The uncertainty at simulator runs is 0
  - unless a nugget is used
- The uncertainty between points decreases as it gets closer to simulator runs
- The width of bounds between points increases as  $\delta$  decreases



- A Gaussian process, with parameters estimated by the training data
- Any point can be estimated
- Uncertainty in the estimate can also be estimated



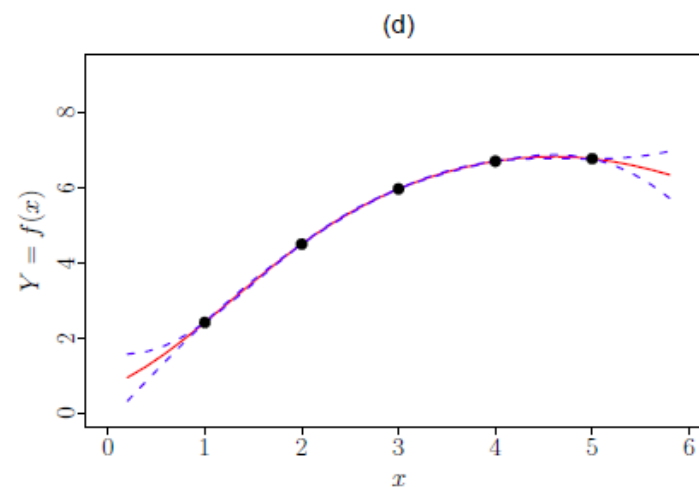
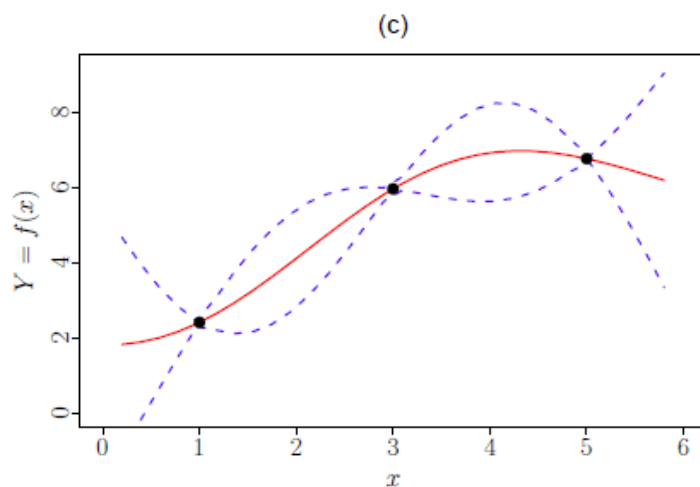
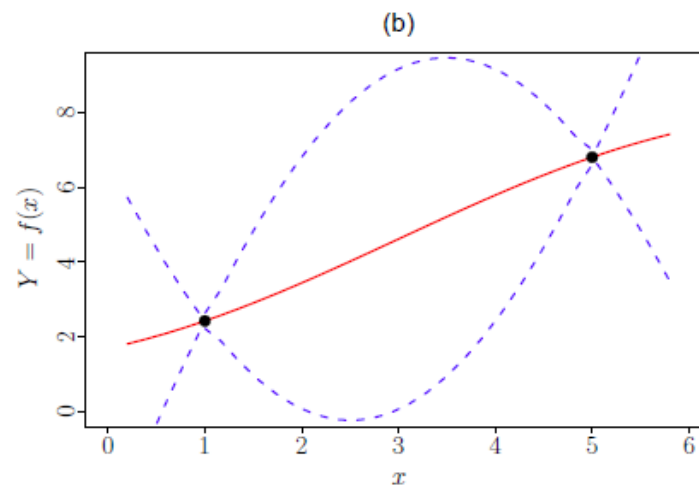
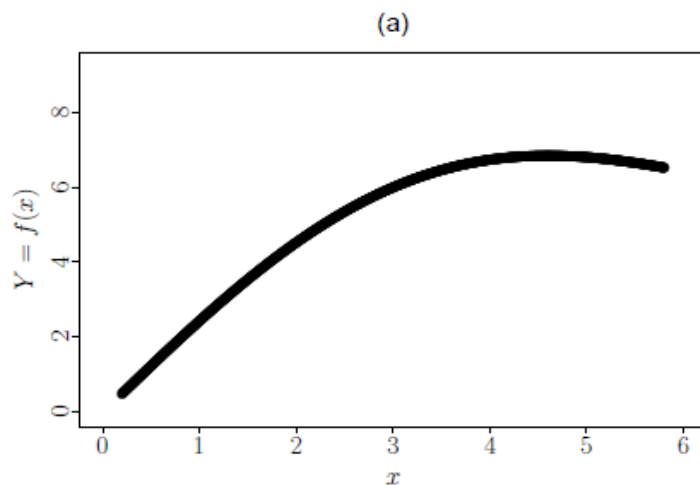
- The emulator is non-unique
- Hyperparameters will be different if you re-run
  - Unless you set the seed
- Using prior information can help stability
- Often, it's not actually a problem...



# Building an emulator



UNIVERSITY OF LEEDS



From Jill Johnson following O'Hagan (2006)

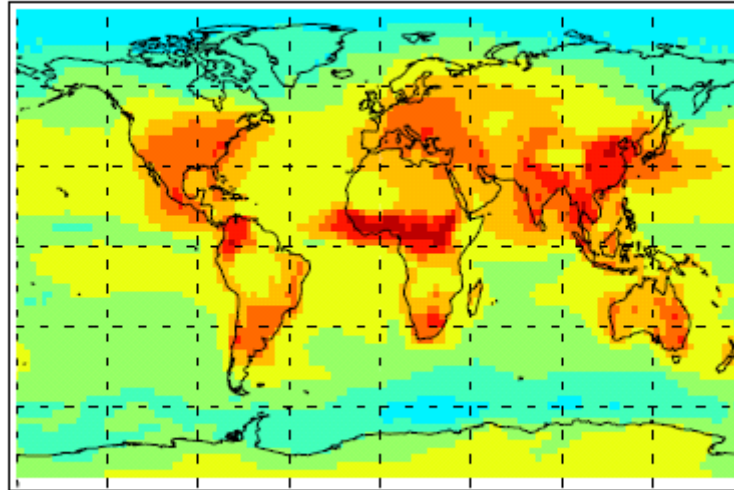
# Aerosol model emulation with DiceKriging (Roustant et al., 2012 & Lee, et al., 2013)



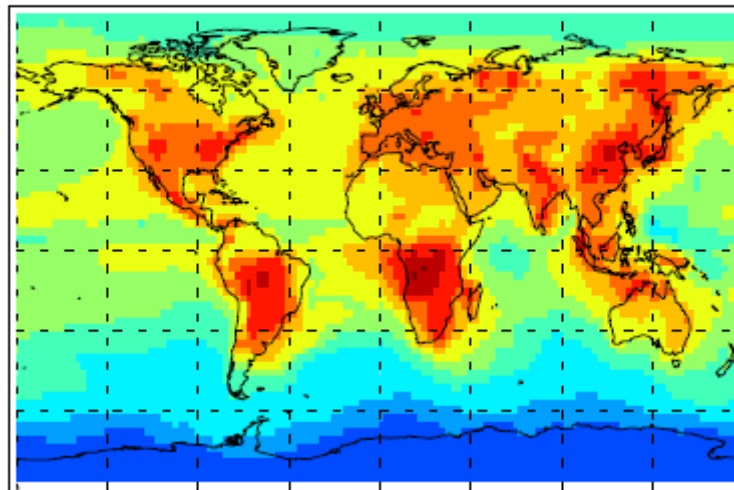
UNIVERSITY OF LEEDS

Emulation done  
gridbox by  
gridbox for  
sensitivity  
analysis

a) Emulator mean CCN ( $\mu_{\text{CCN}}$ )



d) Emulator mean CCN ( $\mu_{\text{CCN}}$ )



- Need to update the prior with simulator points
- Require these points to provide good information about the function through uncertain space
- Often consider no prior information regarding particular parts of the space

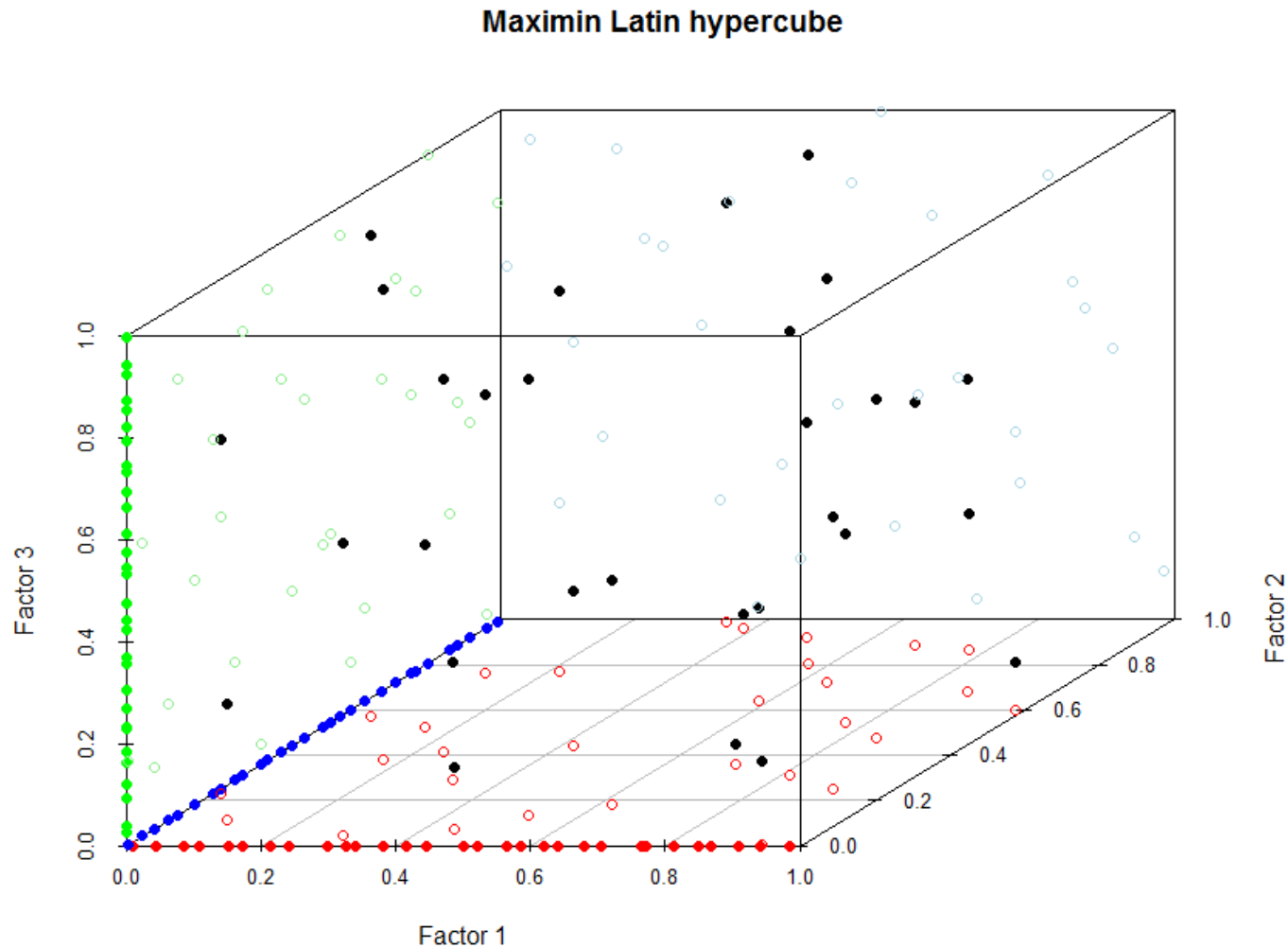


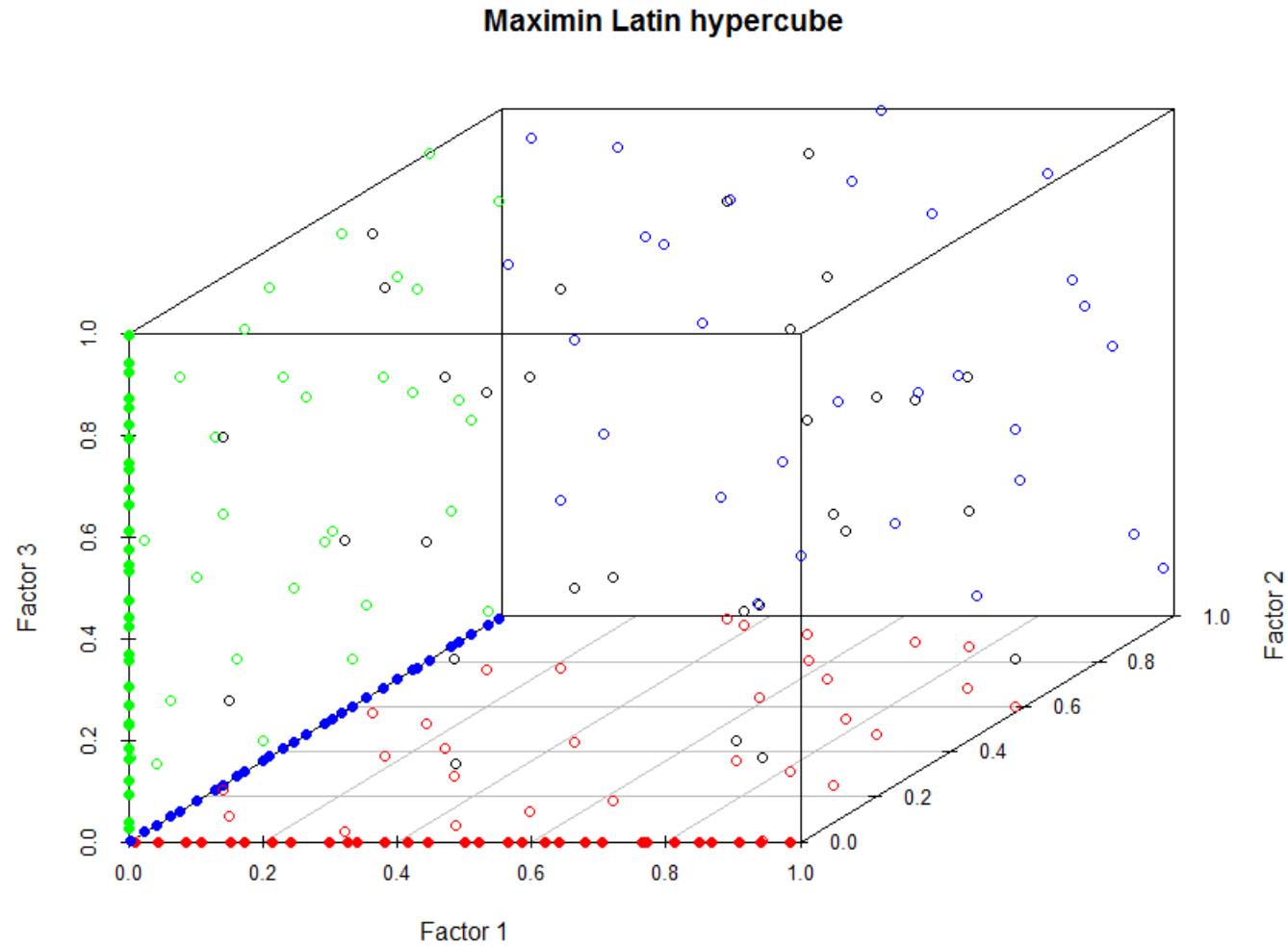


## Latin hypercube sampling

- Good for inactive dimensions
- Use optimum or maximin
- May want piecewise LHS







- When emulator validation suggests can update the design sequentially
- If no reason to search particular region, augment the design
- Can add points sequentially, perhaps Sobol, to particular regions



- These are not recommended
- The covariance function is better estimated using a variety of distances between points
- Uncertainty between training points will be very high



- Have to check the emulator can predict what the emulator would say
- Usually want to check 'out of sample'
- Can use leave one out when circumstances dictate



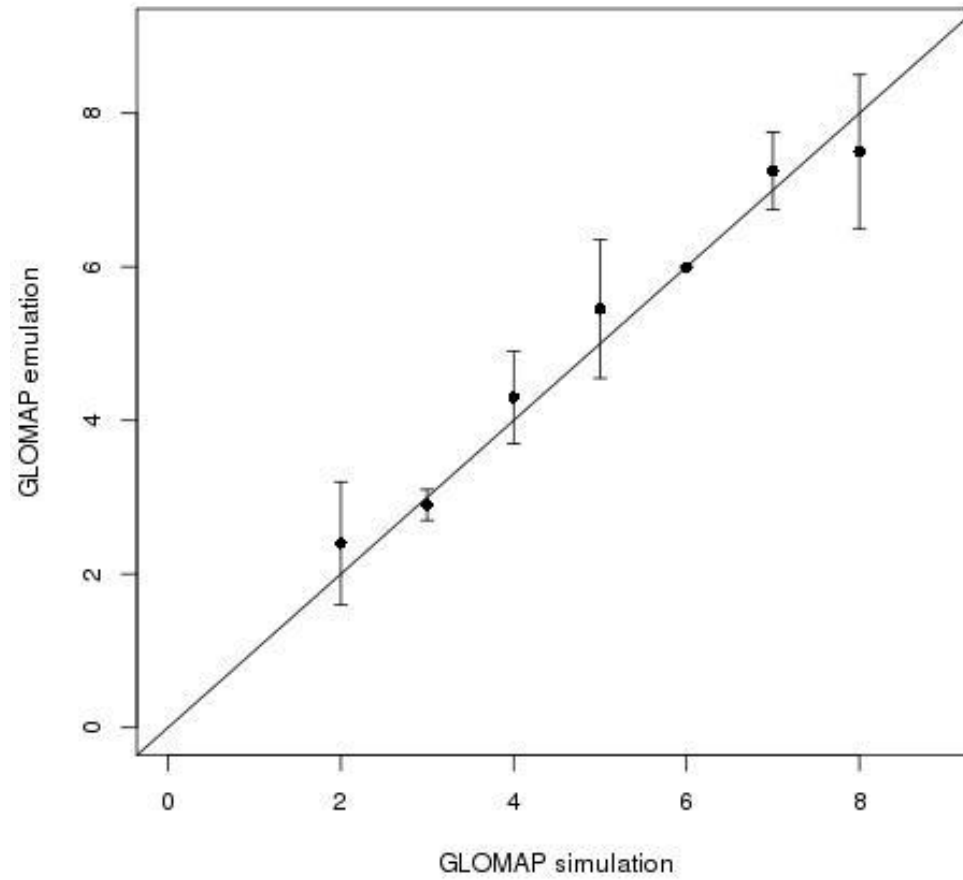
## Prediction value

Ensure emulator  
predictions are  
close to  
simulator

## Prediction uncertainty

Ensure  
uncertainty in  
emulator is small  
enough for  
further inference







- When out-of-sample tend to design two groups, as per Bastos & O'Hagan (2009)
- **1/3 close** to simulated points and **2/3 spaced out**
  - Helps reveal particular difficulties



- Individual prediction errors
  - Treat like regression errors, approximately Gaussian and  $<|2|$  standardised
- Mahalanobis distance
  - single summary of individual prediction errors
  - extreme values indicate emulator/simulator mismatch
- Pivoted Cholesky errors
  - following variance decomposition
  - particular patterns aid interpretation of mismatch



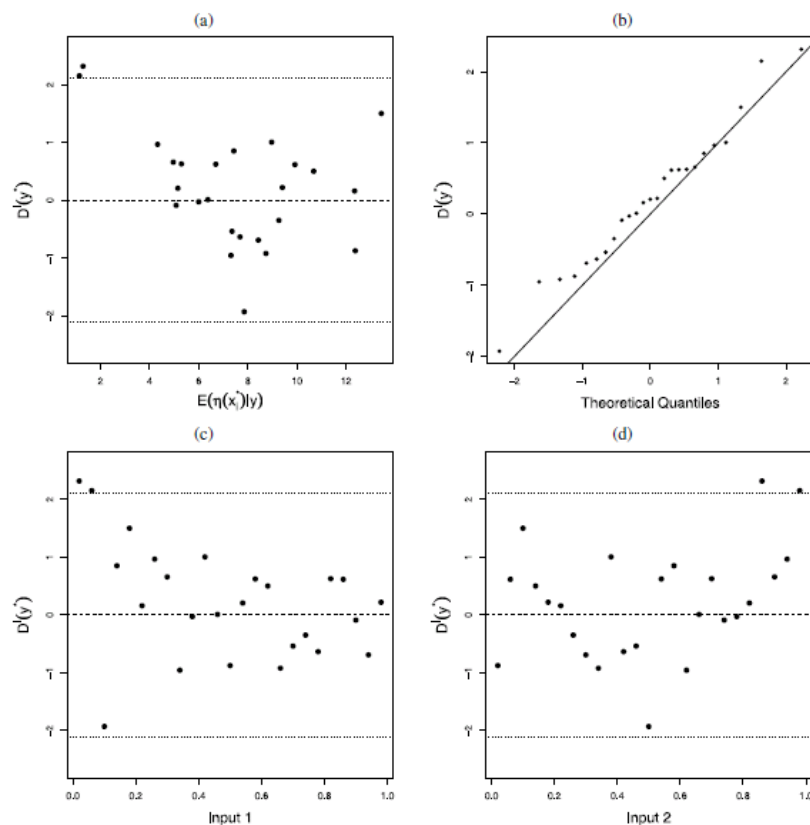
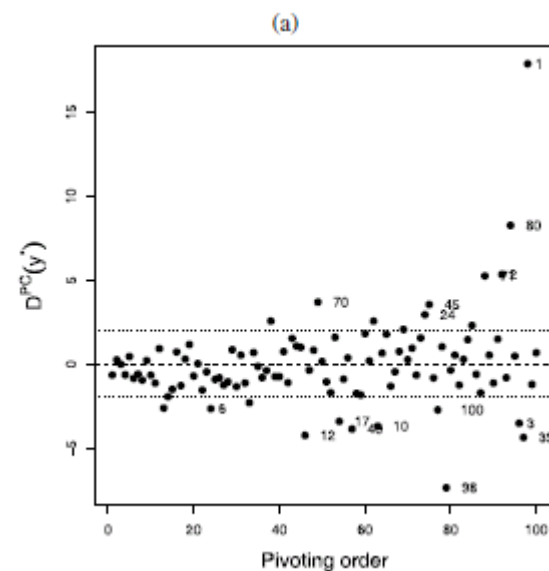
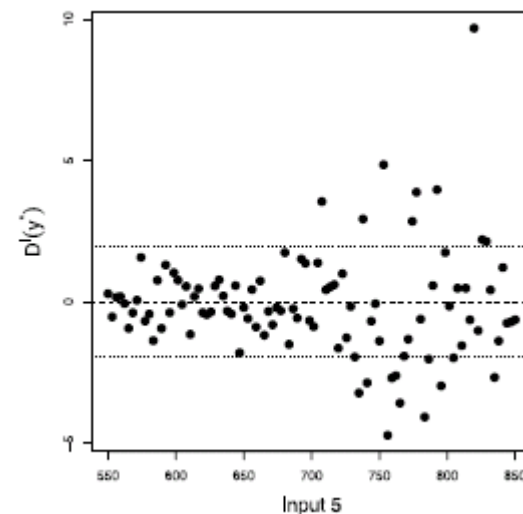


Figure 1. Graphical diagnostics for the toy example using the individual standardized errors: (a)  $D^I(y^*)$  against the emulator predictions; (b) QQ-plot; (c)  $D_I(y^*)$  against input 1; and (d)  $D_I(y^*)$  against input 2.

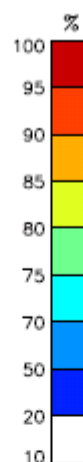
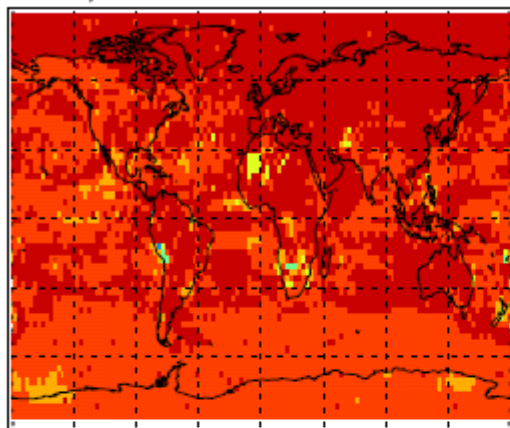


# Aerosol emulator validation

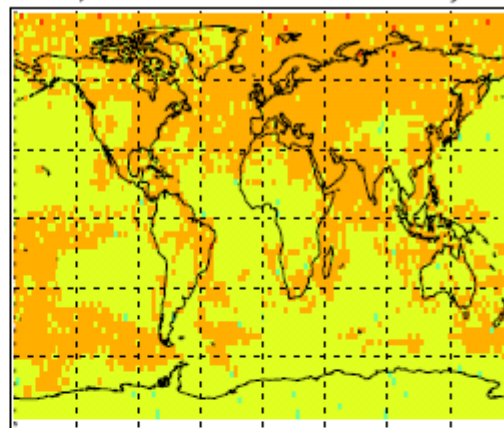


UNIVERSITY OF LEEDS

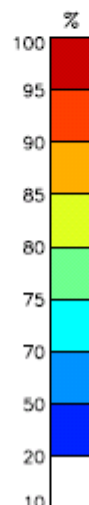
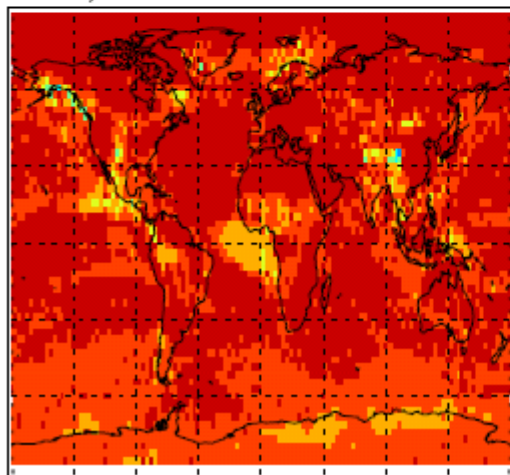
a) JAN emulator validation



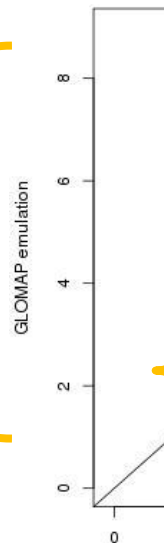
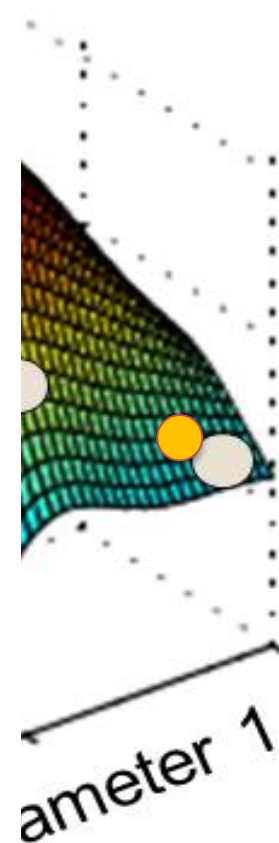
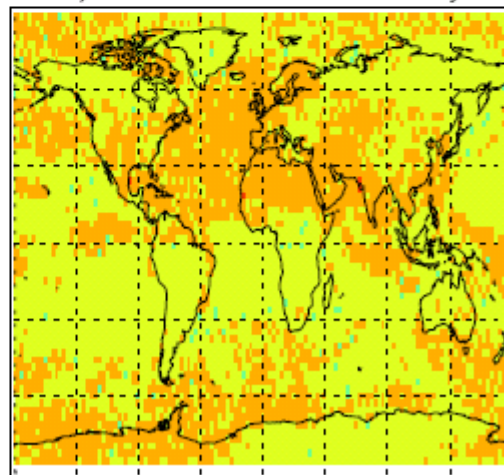
b) JAN emulator uncertainty



c) JUL emulator validation



d) JUL emulator uncertainty



- When out-of-sample is not possible
- Leave a point out, build emulator, produce validation statistics
- Repeat for all points and compare statistics
- Any points with extreme statistics indicate problems

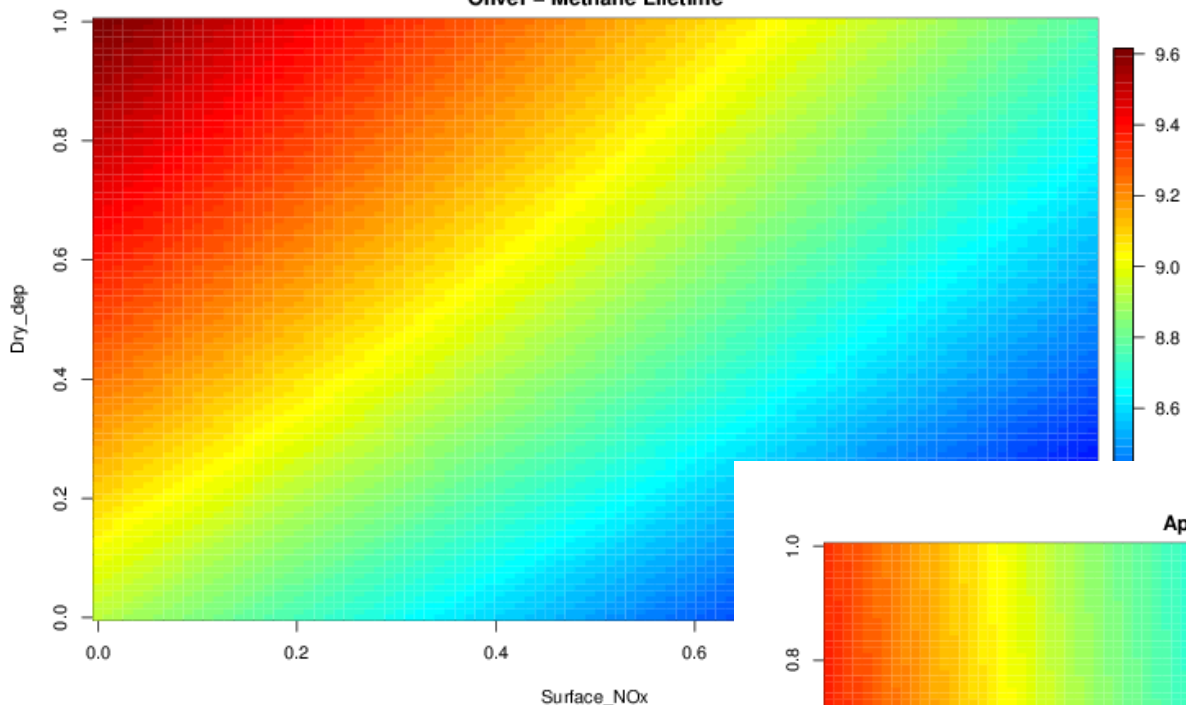


# Emulation for understanding

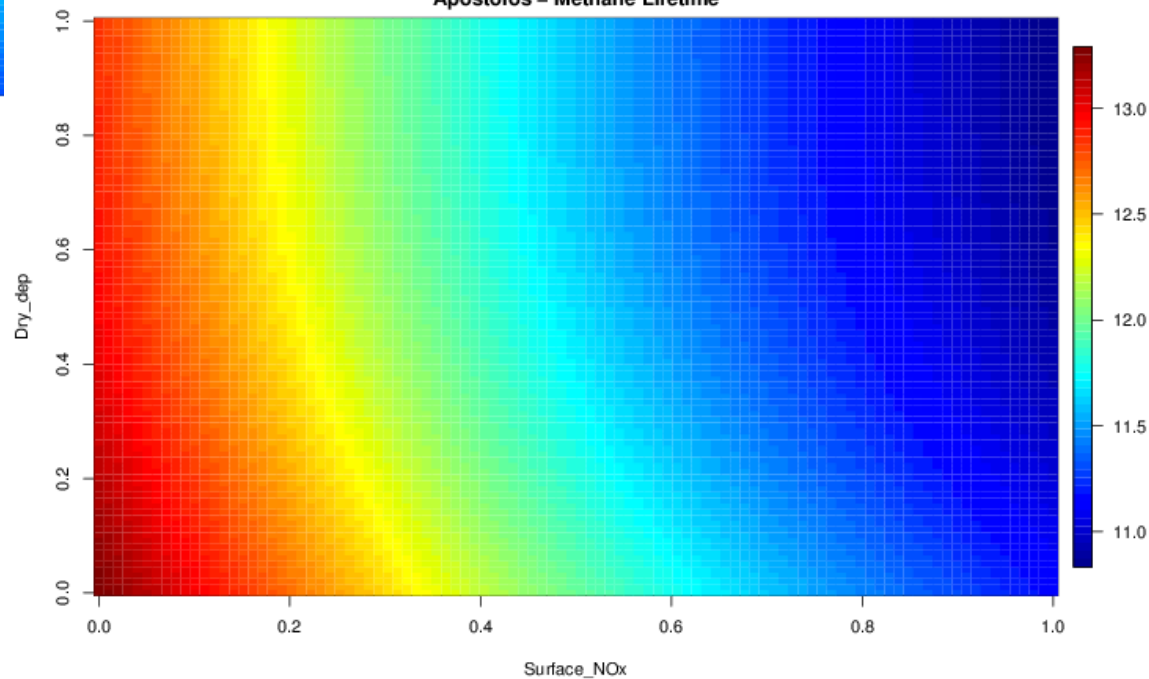


UNIVERSITY OF LEEDS

Oliver – Methane Lifetime



Apostolos – Methane Lifetime

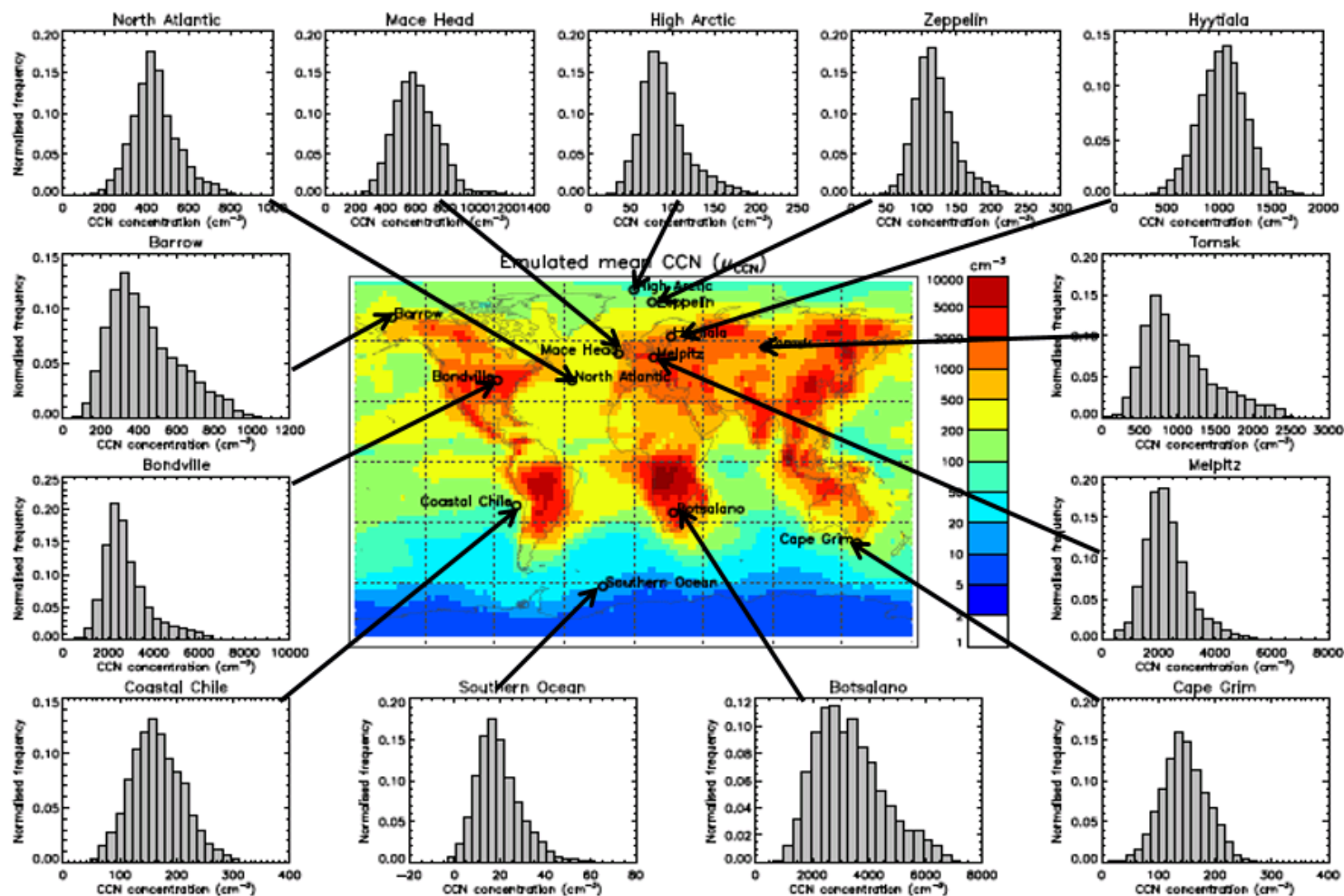


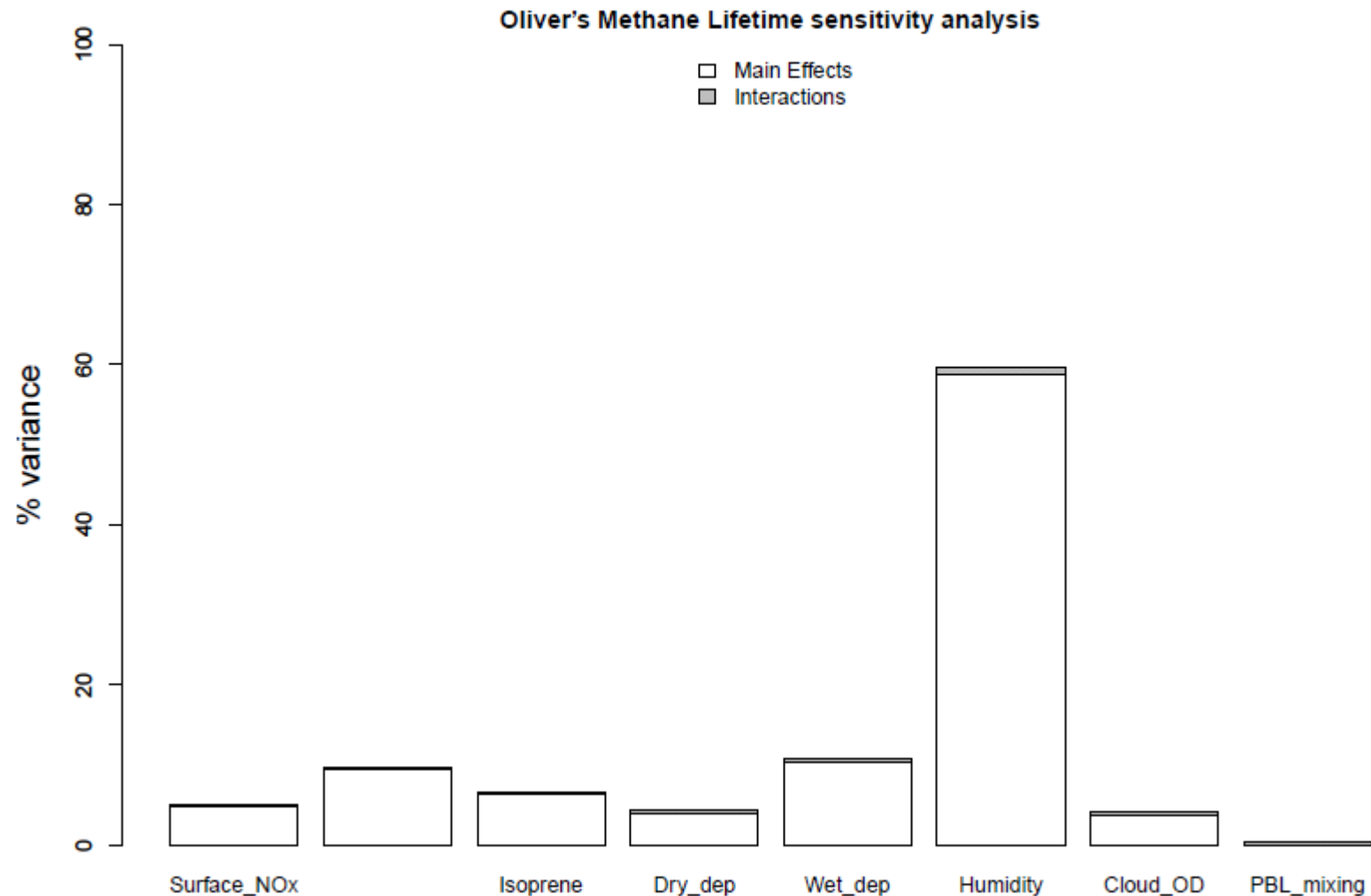
The Leverhulme Trust

# Emulation for uncertainty analysis



UNIVERSITY OF LEEDS





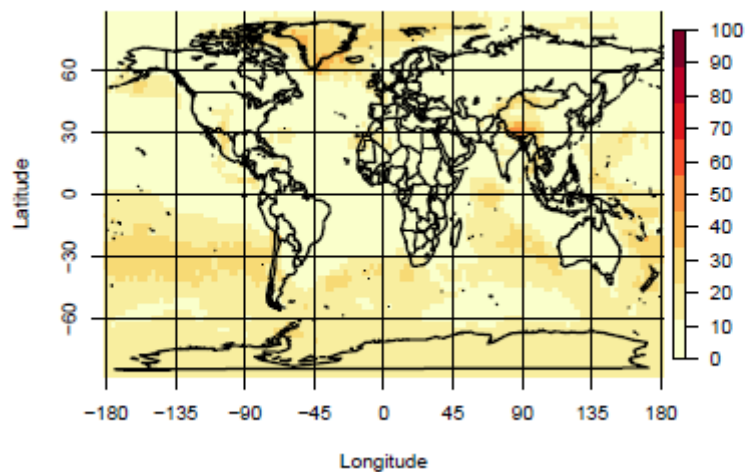


# Aerosol model sensitivity analysis

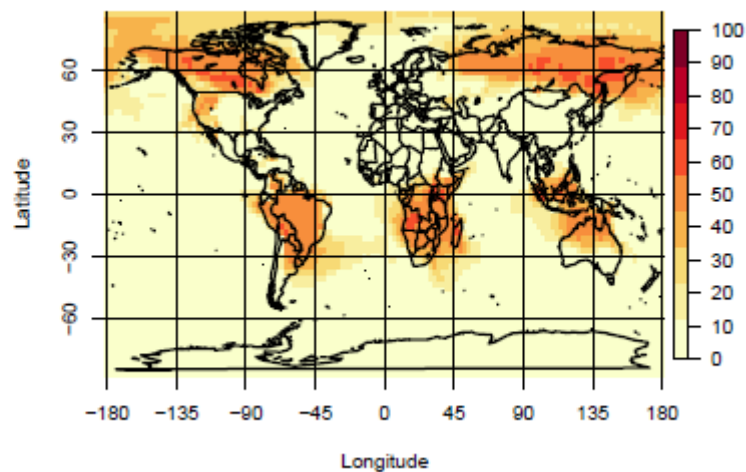


UNIVERSITY OF LEEDS

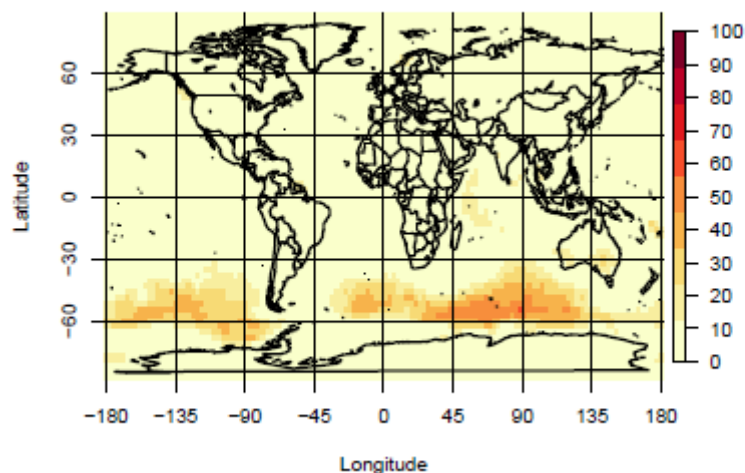
AIT\_WIDTH JULY



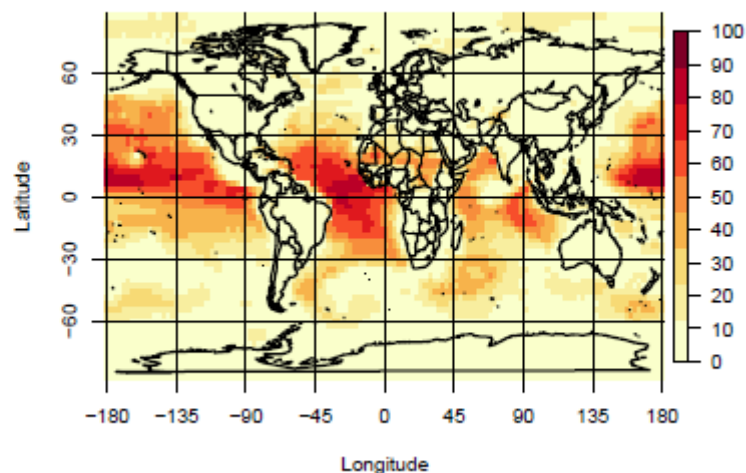
BB\_DIAM JULY



SS\_ACC JULY



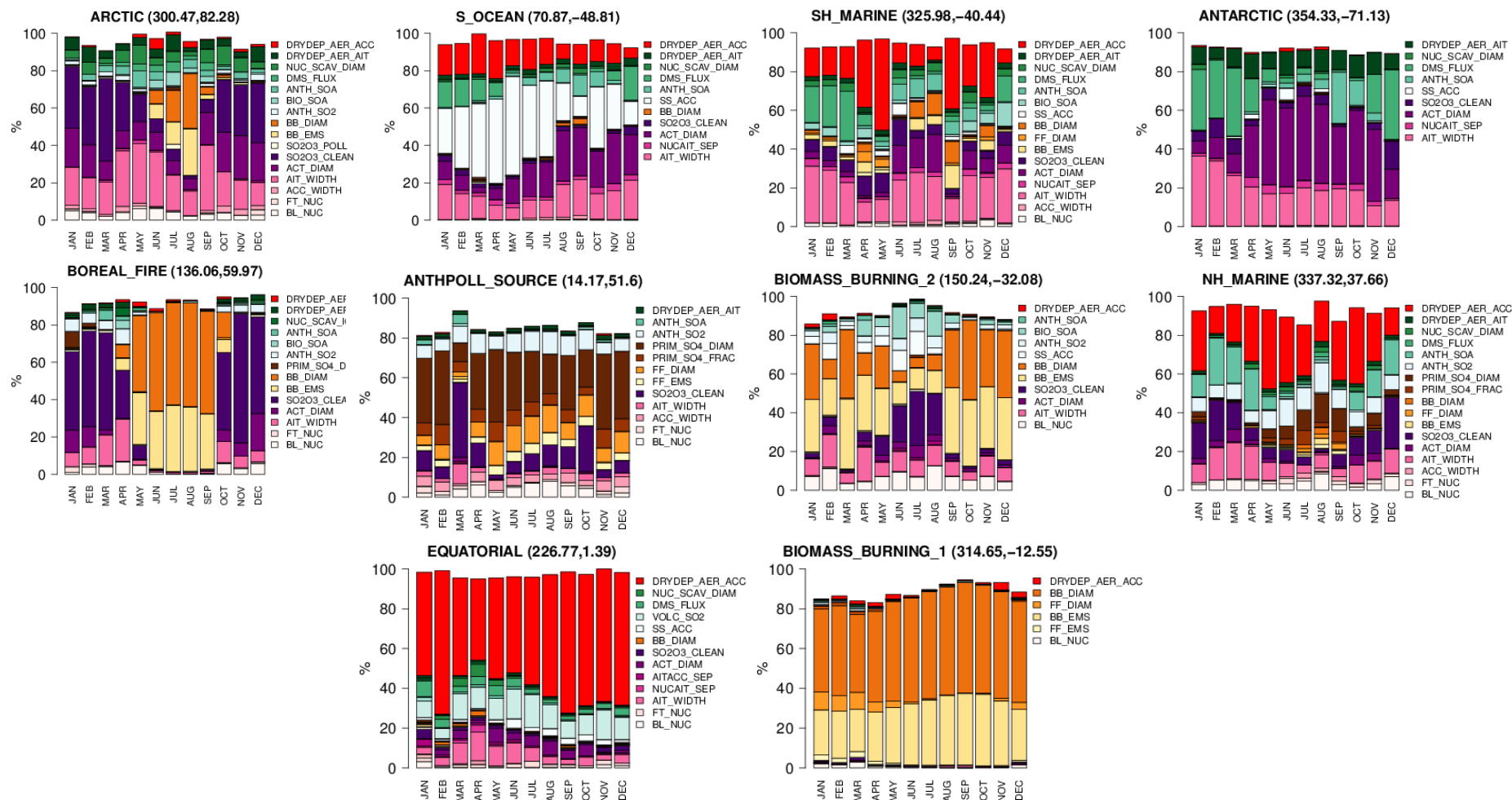
DRYDEP\_AER\_ACC JULY



# Aerosol model sensitivity analysis II



UNIVERSITY OF LEEDS

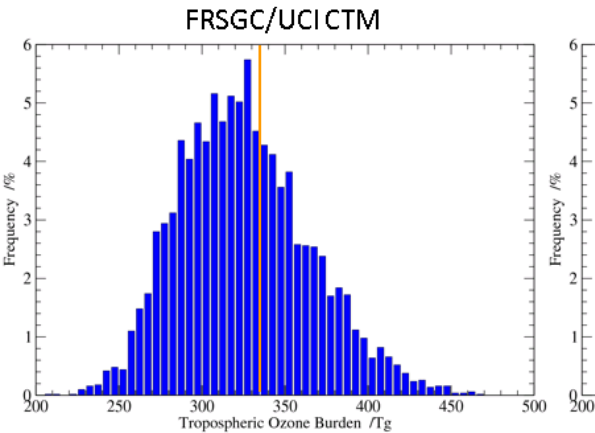


# Multiple model sensitivity analysis I

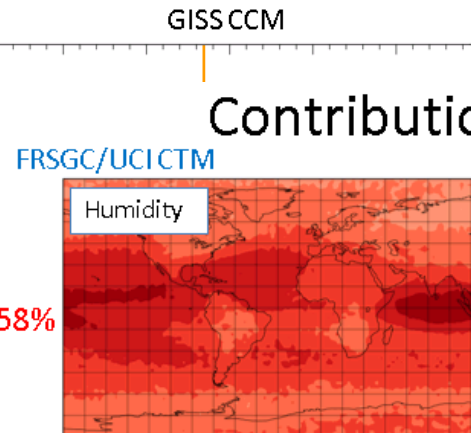


UNIVERSITY OF LEEDS

## Contributions to Uncertainty in $O_3$

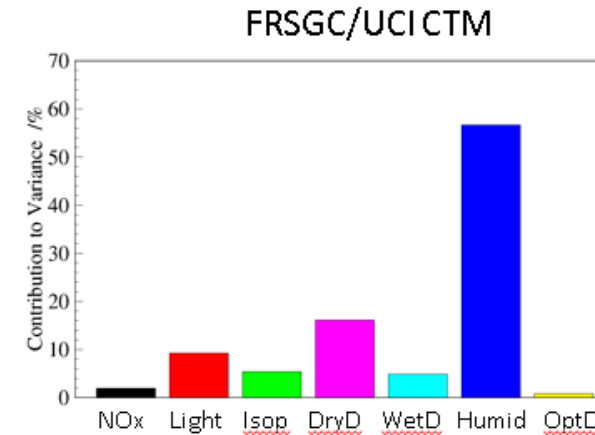
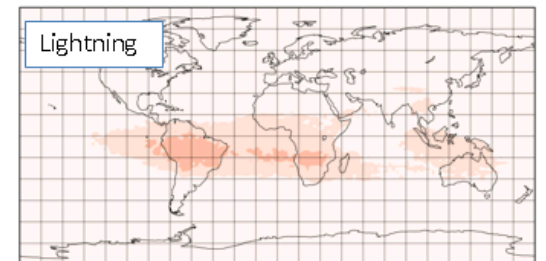
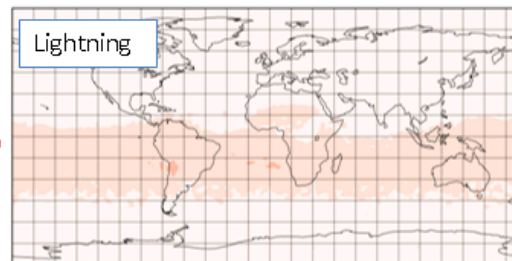
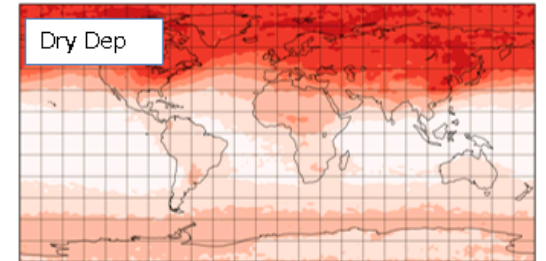
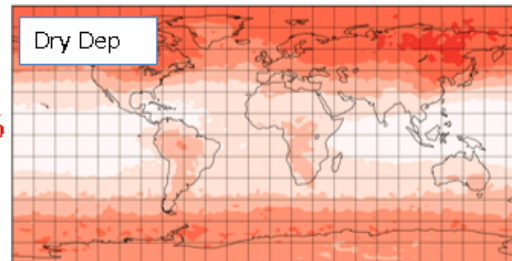
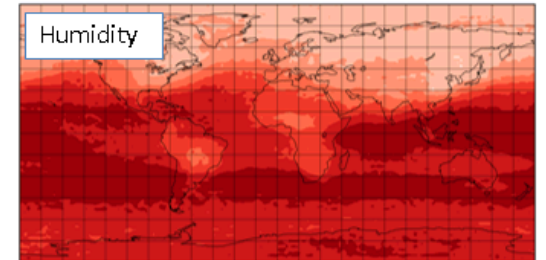
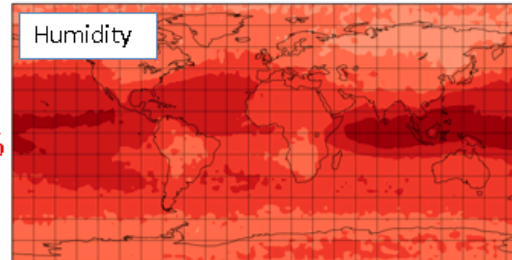


FRSGC:  $325 \pm 39$  Tg



FRSGC/UCI CTM

GISS CCM



# Multiple model sensitivity analysis II

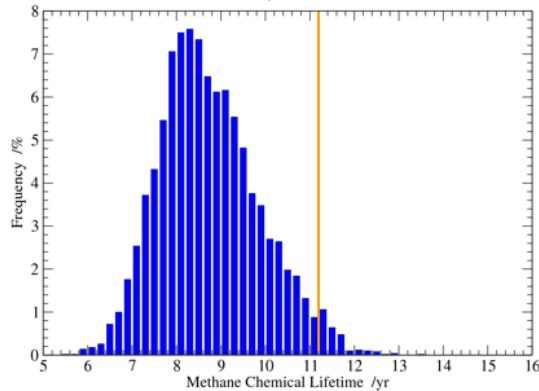


UNIVERSITY OF LEEDS

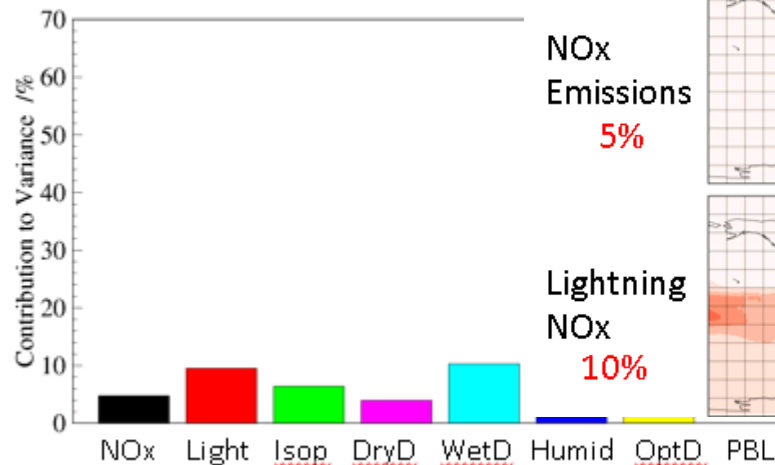
Tropospheric

## Contributions to Uncertainty in CH<sub>4</sub> Loss Rate

FRSGC/UCI CTM



FRSGC/UCI CTM



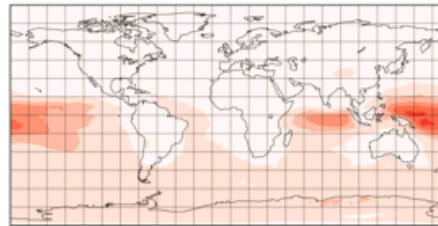
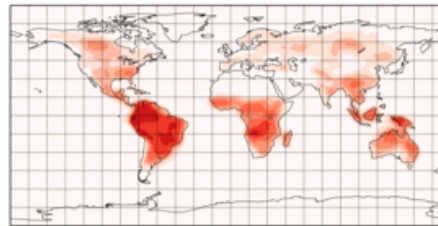
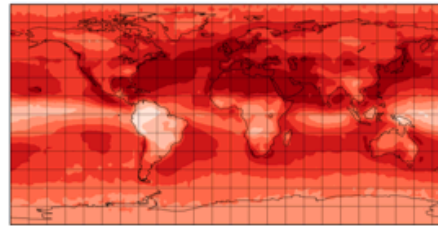
FRSGC/UCI CTM

Humidity  
58%

Isoprene Emissions  
7%

NOx Emissions  
5%

Lightning NOx  
10%



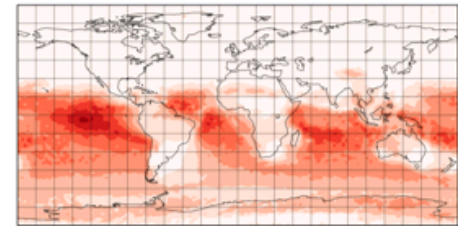
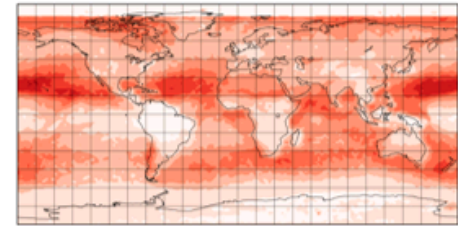
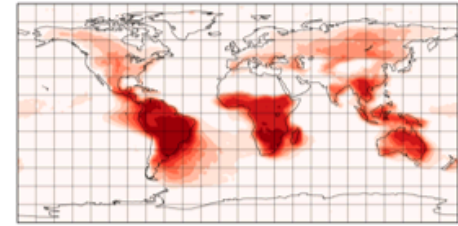
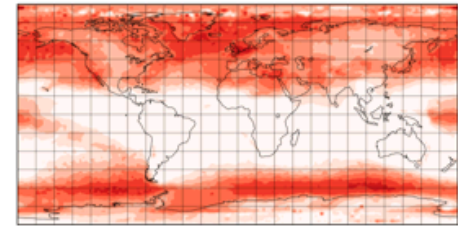
GISS CCM

Humidity  
5%

Isoprene Emissions  
17%

NOx Emissions  
35%

Lightning NOx  
25%



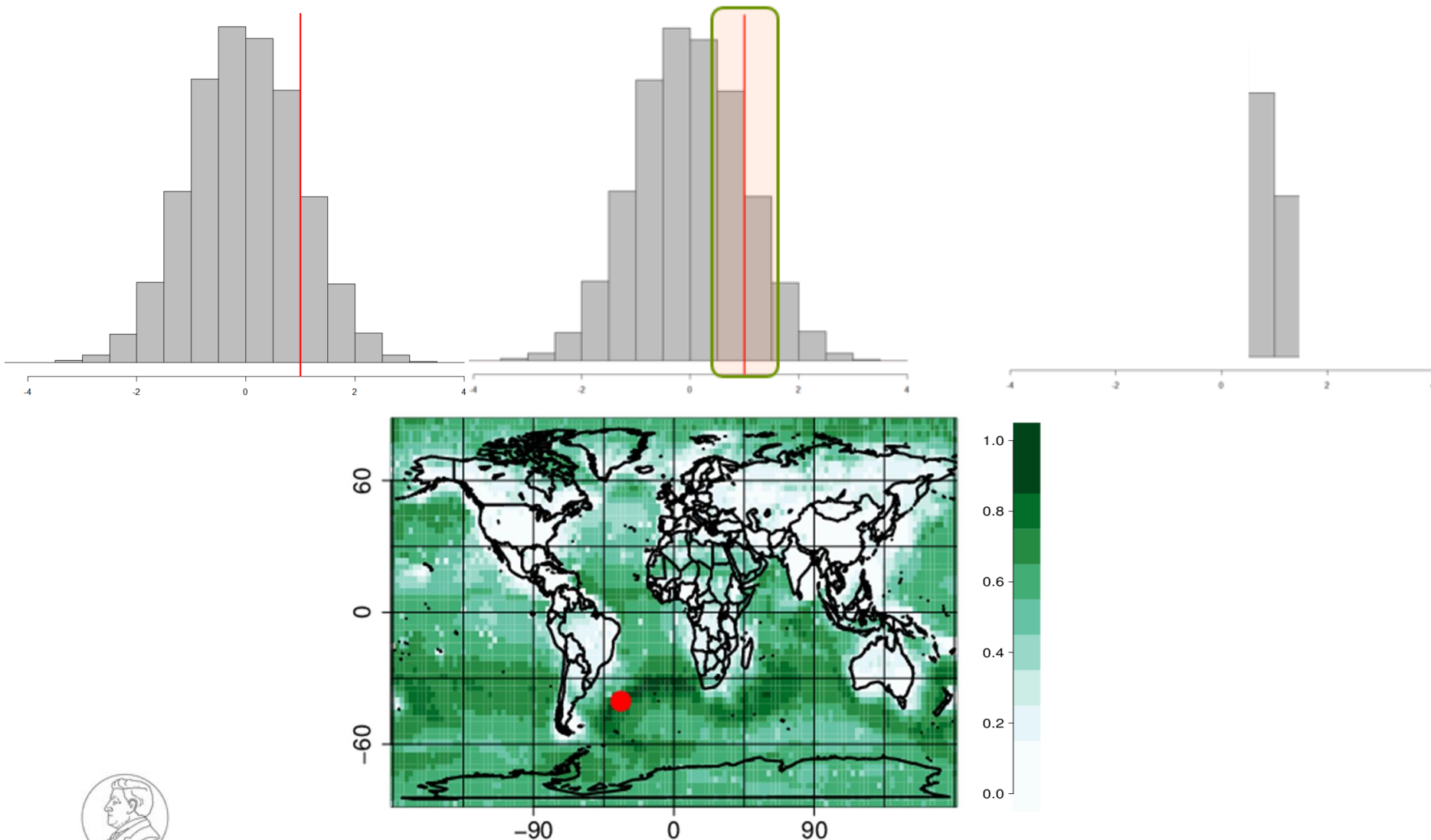
NOx Light Isop DryD WetD Humid OptD PBL



# Emulators for model constraint I



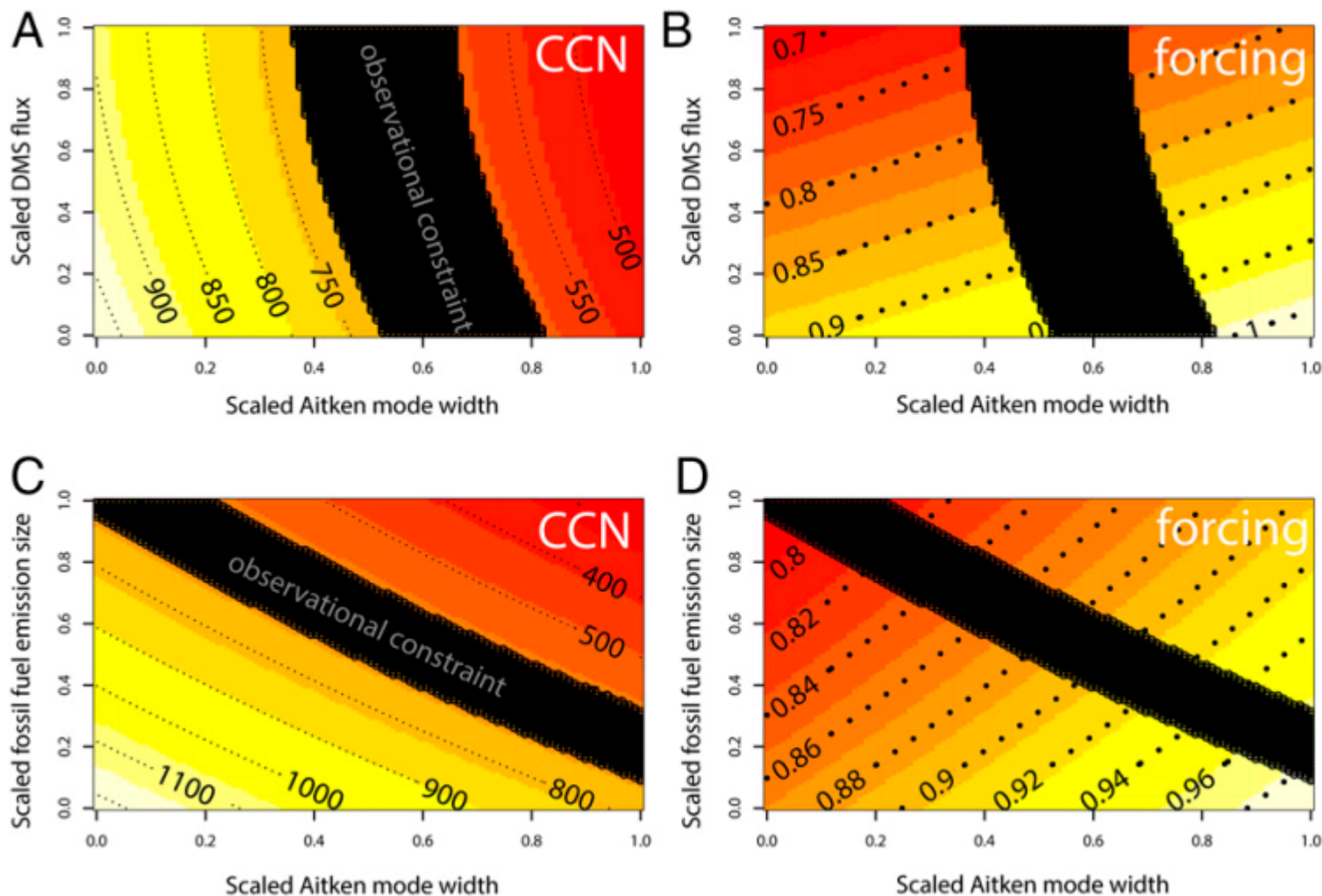
UNIVERSITY OF LEEDS



# Emulators for model constraint II



UNIVERSITY OF LEEDS



- Stochastic simulator – add an estimated nugget to covariance
- Multivariate emulator – must characterise the covariance between multiple outputs
- Discontinuities – could build multiple emulators, perhaps multivariate



- Extrapolate too far with confidence
- Estimate a simulator output it's not trained to
- Replace the simulator





- When the simulator is very cheap to run
- When it won't validate
- When you haven't got a good set of simulator runs for training
- When you aren't sure what you'll use it for



- Rasmussen, Carl Edward, and Christopher KI Williams., 2006. Gaussian processes for machine learning. Vol. 1. Cambridge: MIT press.
- O'Hagan, A., 2006. Bayesian analysis of computer code outputs: a tutorial. Reliability Engineering & System Safety, 91(10), pp.1290-1300.
- Roustant, O., Ginsbourger, D. and Deville, Y., 2012. DiceKriging, DiceOptim: Two R packages for the analysis of computer experiments by kriging-based metamodeling and optimization.
- McKay, M.D., Beckman, R.J. and Conover, W.J., 1979. Comparison of three methods for selecting values of input variables in the analysis of output from a computer code. Technometrics, 21(2), pp.239-245.
- Bastos, L.S. and O'Hagan, A., 2009. Diagnostics for Gaussian process emulators. Technometrics, 51(4), pp.425-438.
- Saltelli, A., Chan, K. and Scott, E.M. eds., 2000. Sensitivity analysis (Vol. 1). New York: Wiley.
- Pujol, G., Iooss, B. and Iooss, M.B., 2008. Package 'sensitivity'.



**Lee LA**; Reddington CL; Carslaw KS (2016) [On the relationship between aerosol model uncertainty and radiative forcing uncertainty](#), *Proceedings of the National Academy of Sciences of the United States of America*, **113**, pp.5820-5827. doi: [10.1073/pnas.1507050113](#)

Carslaw KS; **Lee LA**; Reddington CL; Pringle KJ; Rap A; Forster PM; Mann GW; Spracklen DV; Woodhouse MT; Regayre LA; Pierce JR (2013) [Large contribution of natural aerosols to uncertainty in indirect forcing](#), *Nature*, **503**, pp.67-71.

**Lee LA**; Pringle KJ; Reddington CL; Mann GW; Stier P; Spracklen DV; Pierce JR; Carslaw KS (2013) The magnitude and causes of uncertainty in global model simulations of cloud condensation nuclei, *Atmospheric Chemistry and Physics*, **13**, pp.8879-8914. doi: [10.5194/acp-13-8879-2013](#)

**Lee LA**; Carslaw KS; Pringle KJ; Mann GW (2012) Mapping the uncertainty in global CCN using emulation, *Atmospheric Chemistry and Physics*, **12**, pp.9739-9751. doi: [10.5194/acp-12-9739-2012](#)

**Lee LA**; Carslaw KS; Pringle KJ; Mann GW; Spracklen DV (2011) [Emulation of a complex global aerosol model to quantify sensitivity to uncertain parameters](#), *ATMOSPHERIC CHEMISTRY AND PHYSICS*, **11**, pp.12253-12273. doi: [10.5194/acp-11-12253-2011](#)

