# Correcting a bias in a climate model with an augmented emulator

Doug McNeall[1], Jonny Williams[2], Richard Betts[1,3], Ben Booth[1], and Peter Challenor[3]

[1]Met Office Hadley Centre, FitzRoy Road, Exeter, EX1 3PB, UK
[2]NIWA
[3]University of Exeter

**Correspondence:** Doug McNeall (doug.mcneall@metoffice.gov.uk)

**Abstract.** We develop a method of bias correcting a climate model with a Gaussian process emulator, allowing valid values of input parameters to be found even in the presence of a significant model bias.

A previous study (McNeall et al., 2016) found that a climate model had to be run using land surface input parameter values from very different, almost non-overlapping parts of parameter space in order to satisfactorily simulate the Amazon and other forests respectively. As the forest fraction of the other forests in the model were broadly correct at the default parameter settings and the Amazon too low, that study suggested that the problem likely lay in the model's treatment of the Amazon region than the other regions. The study suggested that this might be due to (1) structural errors such as missing deep-rooting in the Amazon in the land surface component of the model, (2) a warm-dry bias in the Amazon climate of the model, or a combination of both.

In this study we bias correct the climate of the Amazon in a climate model using an "augmented" Gaussian process emulator, where variables often regarded as outputs of the model are treated as inputs. We treat the regional temperature and precipitation of the model as additional inputs to the emulator alongside the standard model inputs. We can then explore the relationship between climate, input parameters and the output of the emulator, forest fraction, finding that the forest fraction is nearly as sensitive to climate variables as any of the land surface inputs. Bias correcting the climate in the Amazon region using the emulator corrects the forest fraction to tolerable levels in the Amazon at many candidates for land surface input parameter values, including the default ones. It also increases the valid input space shared with that suggested by the other forests. We no longer need to invoke a structural model error in the Amazon, beyond having too dry and hot a climate.

Using the augmented emulator allows the bias correction a pre-existing coupled ensemble of climate model runs, reducing the risk of choosing poor parameter values because of an error in a sub-component of the model. We discuss the potential of the augmented emulator to act as a translational layer between model sub-components simplifying the process of model tuning when there are potential compensating errors, and helping model developers prioritise model errors to target. Our technique has the potential to help choose good input parameters for a model, and to efficiently project the impacts of a changing climate, even when there are significant biases in a sub-component of the model.

# 1 Introduction

## 1.1 Uncertainty quantification, compensating errors and model discrepancy

The field of Uncertainty Quantification (UQ) has seen a rapid development of methods to quantify uncertainties when using complex computer models to simulate real, physical systems. These models often contain simplifications of processes too complex to represent explicitly in the model, termed parameterisations. Associated with these parameterisations are coefficients called input parameters, the values of which are uncertain and can be set by the model developer. The settings have a material effect on the way the parameterisations operate, and therefore on output of the model, but often to an extent that is unknown until the model is run. Input parameters are subject to uncertainty and may be difficult or even impossible to observe, having no direct analogue in the real system. The process of setting the values of the input parameters so that the simulator output best matches the real system is called tuning, and where a probability distribution is assigned for the input parameters, it is termed calibration. Uncertainty in input parameters can induce uncertainty in the output of the model, leading to uncertainty in projections of future climate states or reconstructions of past ones.

Without strong prior information it can be difficult to attribute a difference between simulator output and the real system to underlying model errors, to an incorrect set of input parameters, or to inaccuracies in the observations. Similarly, there are often a number of ways to set parameters that lead to a particular model output, with poor choices of a particular parameter compensating for poor choices of other parameters, or for modelling errors. This situation means that a good candidate for input parameters might be found in a large volume of input space, and projections of the model made with candidates from across that space might display a very wide range of outcomes. This problem is sometimes referred to as "identifiability", but otherwise known as "equifinality", or the "degeneracy" of model error and parameter uncertainty. It can be relatively easy to find a good subset of input parameters given a small set of inputs and outputs and a well behaved relationship between the two. This situation might be found for a subcomponent of a climate model, where there are good observations of the system being studied, for example. Improving a coupled climate model however can require an involved and lengthy process of development. Some components of the model may have been tuned to compensate for errors in others or there may be unknown errors in the model or observations. Further, more complex models are computationally expensive and so infeasible to run in enough configurations to be able to identify these errors.

Hourdin et al. (2017) offer a summary of current practice in the somewhat understudied and sparsely documented field of climate model tuning. While there are clearly common features, there appear no standard procedures for climate model tuning however. As Hourdin et al point out, it remains an art as well as a science. Various individual centres have begun to document their tuning practices with regard to tuning targets and procedures (Schmidt et al., 2017; Zhao et al., 2018; Walters et al., 2017).

Parameter tuning occurs at different stages in model development, perhaps starting with single column version of the model. The climate model components to be coupled might be then tuned with standard boundary conditions - for example tuning a land/atmosphere component with fixed or historically observed sea surface temperatures. Finally, a system-wide tuning might be used to check that there are minimal problems once everything has been coupled together.

Golaz et al. (2013) Show the potential impact of compensating errors in tuning. They find that two different but plausible parameter configurations of the cloud formations of the coupled climate model GFDL-CM3 can result in similar present-day radiation balance. The configurations did not differ in their present day climate, but showed significantly different responses to historical forcing and therefore historical climate trajectories.

5    Although climate model tuning is overall a subjective process, individual parts are amenable to more algorithmic approaches. Statistical and machine learning approaches to choosing parameters to minimise modelling error, or to calculate probability distributions for parameters and model output are known as uncertainty quantification (UQ).

The problem of accounting for model discrepancy when using data to learn about input parameters is becoming more widely recognised in UQ. It was formalised in a Bayesian setting by Kennedy and O'Hagan (2001). The authors suggested

10   simultaneously estimating a model discrepancy - there called model inadequacy - as a function of the inputs, using a Gaussian process prior.

Arendt et al. (2012a) offer a number of examples of identifiability problems, ranging from solvable using mild assumptions through to virtually impossible. In a companion paper (Arendt et al., 2012b), they outline a way of improving identifiability using multiple model responses.

15   Brynjarsdóttir and O'Hagan (2014) argued that only by accounting for model discrepancy does even a very simple simulator have a chance of making accurate predictions. Further, they found that only where there is strong prior evidence about the nature of that model discrepancy is it possible to solve the inverse problem and recover the correct inputs. Without this strong prior evidence the estimate of the correct parameters is likely to be overconfident, and wrong, leading to overconfident and wrong predictions of out-of-sample data.

20   Some of the dangers of overconfident and wrong estimates of input parameters and model discrepancy can be reduced using a technique called history matching (Craig et al., 1996), sometimes called pre-calibration or iterated refocussing. The aim of history matching is not to find the most likely inputs, but to reject those unlikely to produce simulations statistically close to observations of the real system. An implausibility measure (I) is calculated, taking into account the distance between the simulator output and the observation, but allowing for uncertainty in the observations, the simulator output and the simulator

25   discrepancy. Those inputs that produce a large implausibility score are ruled out from consideration as candidate points.

An excellent introduction and case studies can be found in Andrianakis et al. (2015), or in Vernon et al. (2010) History matching is perhaps less ambitious but correspondingly more robust than calibration methods, and a full calibration can be carried out once the history matching procedure has been completed.

(McNeall et al., 2013) studied an ensemble of an ice sheet model and found that using a single type of observation for ruling

30   out input space was not very powerful - particularly if there was not a very strong relationship between an input parameter and the simulator output. A key technique therefore is to use multiple data sets for the history matching, ruling out a candidate input space according to an empirical rule. Several rules have been used - for example using the maximum implausibility of a multiple comparison, a candidate input point point may be ruled out by a single observation. A more conservative approach is to use the second or third implausibility score, or to use a multivariate implausibility score, both introduced in Vernon et al.

35   (2010). The aim of these scores is to ensure that an unidentified model discrepancy does not result in ruling out candidate points

**3**

that are in fact perfectly good. History matching can be effective in reducing the volume of parameter space that is considered plausible to produce model runs that match the real system. For example, Williamson et al. (2015) report very large reductions (around 99%) in the volume of space considered plausible, when history matching is used in an iterated fashion.

While history matching has often been used used to explore and reduce the input parameter space of expensive simulators, its use as a tool to find discrepancies, bias and inadequacies in simulators is less developed. Williamson et al. (2015) argue that what was assumed a structural bias in ocean model HadCM3 could be corrected by choosing different parameters. In a different system McNeall et al. (2016) argue that a standard set of parameters for the land surface component of the climate model FAMOUS should be retained, and that a bias seen in the simulation of the Amazon rainforest is a simulator discrepancy not a poor parameter choice.

In that case, the model simulated other forests at the standard set of parameters well, and only a tiny volume of parameter space could be found that (barely) adequately simulated all the forests. When cast as a choice between keeping the default parameters, or rejecting them and accepting the new region of parameter space, they argued that the former was more likely to produce a good model, as presumably scientific judgement and expertise informed the original choice of parameters, whereas there were a number of reasons one might reject the proposed parameter space.

## 1.2 Aims of the paper

A well simulated and vigorous Amazon forest at the end of the spinup phase of a simulation experiment is a prerequisite for using the model to make robust projections of future changes in the forest. The analysis of McNeall et al. (2016) (hereafter M16) identified that the land surface input spaces where FAMOUS forest fraction was consistent with observations were very different in the Amazon than they were for other forests. The area of overlap of these spaces - one that would normally be chosen in a history matching exercise - did not simulate any of the forests well, and did not contain the default parameters. M16 suggested that assuming an error in the simulation of the Amazon forest would be a parsimonious choice. Two obvious candidates for the source of the discrepancy in the Amazon were identified: (1) a lack of deep rooting in the Amazon, meaning that trees could not access water at depth as in the real Amazon and (2) a bias in the climate of the model, impacting the vigour of the trees.

This paper revisits and extends the analysis of M16 to attempt to simultaneously (1) assess the impact of a bias corrected corrected climate on the Amazon forest and (2) to identify regions of input parameter space that should be classified as plausible, given a corrected Amazon climate. To bias correct the climate we develop a new method to augment a Gaussian process emulator, with simulator outputs acting as inputs to the emulator alongside the standard input parameters. We use simulated output of forests at different geographical locations to train the emulator, describing a single relationship between the climate of the simulator, the land surface inputs and the forest fraction. In doing so, we develop a technique that might be used to bias correct existing ensembles of coupled models, allowing a more computationally efficient method for final system-tuning of models.

In section **??**, we review the literature on the possible causes of the low Amazon forest fraction in FAMOUS. In section **??**, we describe how we use the temperature and precipitation to augment the Gaussian process emulator. In section **??** we use the

4

emulator to estimate the sensitivity of forest fraction to changes in land surface and climate parameters. In section **??** we use the augmented emulator to bias correct the climates of the forest and examine the effect of that bias correction on the input space that is deemed statistically acceptable in a history matching exercise. In section **??** we search for regions of parameter space where the bias corrected simulator might perform better than at the default parameters. In section **??** we look at regions of climate space where the default parameters would produce statistically acceptable forests. Finally, we offer some discussion of our results in section **??** and conclusions in section **??**.
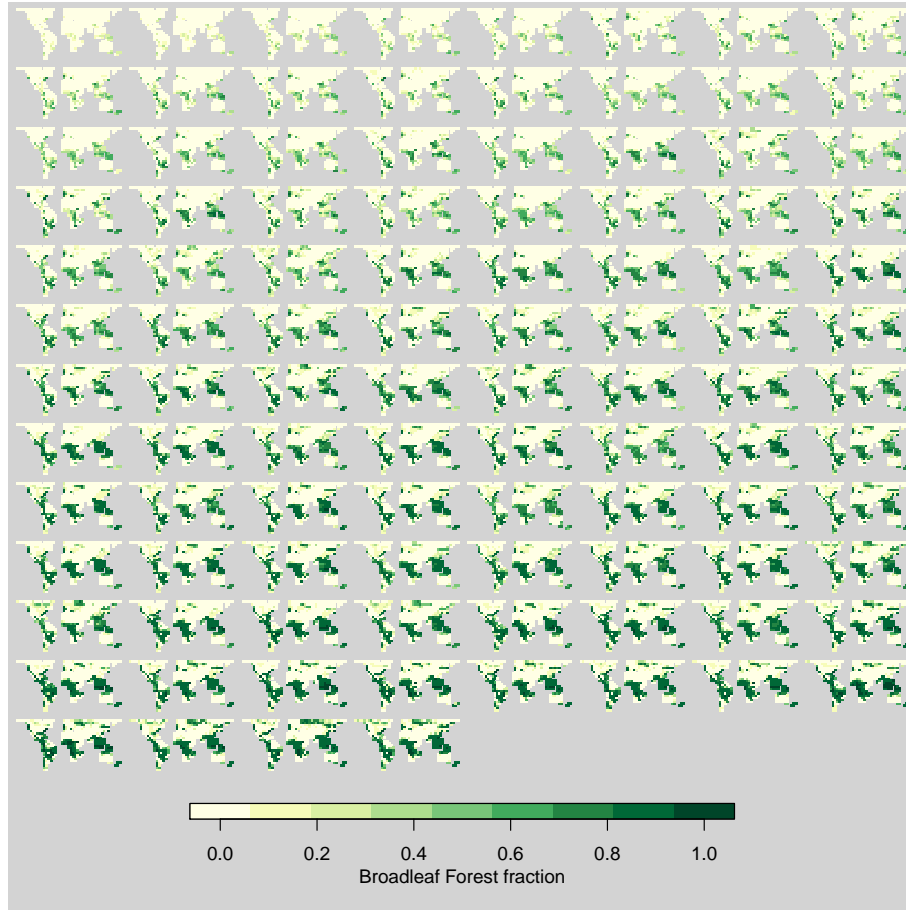
## 2  Climate and forest fraction in FAMOUS

Previous studies have concluded that the climate state has an influence on the Amazon rainforest. Much of that work has been motivated by the apparent risk of dieback of the Amazon forest posed by a changing climate [e.g. Malhi 2008, Cox et al 2004]. We assume here that factors that might precipitate a future change in the simulated state of the Amazon rainforest might also influence the simulation of the steady-state preindustrial forest in FAMOUS. Parameter perturbations and CO2 concentrations have also been shown to influence the simulation of tropical forests in climate models, with increases in CO2 fertilisation and associated increased water use efficiency through stomatal closure offsetting the negative impacts of purely climatic changes. A metric linked to rainforest sustainability by Mahli (2009b) is Maximum Cumulative Water Deficit, which describes the most negative value of climatological water deficit measured over a year. In a similar vein Good et al. (2011, 2013) find that in Hadley Centre models, sustainable forest is linked to dry-season length, a metric which encompasses both precipitation and temperature, along with sensitivity to increasing CO2 levels. No forest is found in regions that are too warm or too dry, and there is a fairly distinct boundary between a sustainable and non-sustainable forest. Galbraith et al (2010) found that temperature, precipitation and humidity had greatly varying influences, and by different mechanisms on changes in vegetation carbon in the Amazon across a number of models, but that rising CO2 mitigated losses in biomass. Poulter et al. (2010) found that the response of the Amazon forest to climate change in the land surface model LPJml was sensitive to perturbations in parameters, but that the dynamics of a dieback in the rainforest was robust across those perturbations. In that case, the main source of uncertainty of dieback was uncertainty in climate scenario. Boulton et al. (2017) found that temperature threshold and leaf area index parameters both have an impact on the forest sustainability under projections of climate change in the Earth system version of HadCM3.

### 2.1  Biases in FAMOUS

M16 speculated that both local climate biases and missing or incorrect processes in the land surface model - such as missing deep rooting in the Amazon - might be the cause of the simulated low forest fraction in the Amazon region at the end of the pre-industrial period in an ensemble of the climate model FAMOUS. In this study we use the ensemble of FAMOUS previously used in M16, to attempt to find and correct the cause of persistent low forest fraction in the amazon, identified in that paper.

The Fast Met Office UK Universities Simulator, FAMOUS (Jones et al., 2005; Smith et al., 2008), is a reduced-resolution climate simulator based on the climate model HadCM3 (Gordon et al., 2000; Pope et al., 2000). The model has many features

**Figure 1.** Broadleaf forest fraction in the FAMOUS ensemble, ranked from the smallest to largest global mean value.

of modern climate simulators, but is of sufficiently low resolution to provide fast and simple data sets with which to develop UQ methods. Full details of the ensemble can be found in M16 and Williams et al. 2013.

The ensemble of 100 members perturbed 7 land surface and vegetation inputs, which had a strong impact on vegetation cover at the end of a spinup period, with atmospheric $CO_2$ at preindustrial conditions (figure). The broadleaf forest fraction in individual ensemble members varies from almost non-existent to vigourous. The strong relationships between forest fraction in each forest and global values implies that perturbations in input parameters exert a larger control over all forests simultaneously, and individual forests to a smaller extent.

M16 extracted aggregated forest fraction data for the Amazon, Southeast Asian, North American and central African forests, along with the global mean. They were only able to find very few land surface parameter settings which the emulator suggested should lead to an adequate simulations of the Amazon forests and the other forests together. Further, these parameter sets were at the edges of sampled parameter space, where larger uncertainty in the emulator may have been driving the acceptance of the parameter sets.

**6**

The ensemble did however have a further perturbation - a parameter denoted "beta", which indexed into one of ten of the best-performing atmospheric parameter sets used in a previous ensemble with the same model. The beta parameter then summarised perturbations in a number of other parameters that impacted the climate of the model. Variations in the parameter did not correlate with any of the land surface parameters in the ensemble, and so was excluded from the analysis in M16.

In this study, we use the same ensemble of forest fraction data used in M16. However, we add temperature and precipitation data, present in the original ensemble but not used to build an emulator in the M16 study, to further our understanding of the causes of the low forest fraction in the Amazon region. The temperature and precipitation data summarise the effects of the parameter on the atmospheric component of the model, in a way that is directly seen by the land surface component of the model. We consider only regions dominated by tropical broadleaf forest, so as not to confound analysis by including other forests which may have a different set of responses to perturbations in parameters, rainfall and temperature.

For temperature observations we use the CRU global monthly surface temperature climatology (Jones 1999), covering the years 1960-1990 . For precipitation we use the average monthly rate of precipitation, covering the years from 1979-2001 from GPCP Version 2.2 provided by the NOAA/OAR/ESRL PSD, Boulder, Colorado, USA, from their Web site at https://www.esrl.noaa.gov/psd/ (Adler et al., 2003). Vegetation fraction observations are adapted from Loveland et al. (2000), and are shown in figure 2. Although the observations all cover slightly different time periods, we expect the differences caused by harmonising the time periods to be very small compared to other uncertainties in our analysis, and to be well covered by our uncertainty estimates.

Figure 2 Observations of broadleaf forest fraction on their native grid (top), and regridded to the FAMOUS grid (bottom).

A plot of regional mean temperature and precipitation in the tropical forest regions in the FAMOUS ensemble (figure 3) indicates the form of the impact that the regional climate has on forest fraction in the climate model. Central African and Southeast Asian climates in the model simulations run in a sweep across the middle of the plot, from dry and cool to wet and warm.

It appears that a wetter climate - which would be expected to stabilize forests - broadly compensates for the forest reductions induced by a warmer climate. Within the ensemble of Central African forests for example, forest fraction increases towards the "cooler, wetter" (top left) part of the climate phase space. Beyond a certain value however, there are no simulated forests in this climatic region. It is clear from the plot that while central African and Southeast Asian forests are simulated in the large part considerably warmer than recent observations, they are also simulated considerably wetter, which might be expected to compensate forest stability. In contrast, while simulated considerably warmer, the Amazon is also slightly drier than recent observations, which might further reduce forest stability.

## 3   Methods

The climate model FAMOUS is computationally expensive enough that we cannot run it for a large enough number of input parameter combinations to adequately explore parameter space and find model biases. To increase computational efficiency we build a Gaussian process emulator: a statistical function that predicts the output of the model at any input, with a corresponding

estimate of uncertainty. The emulator models climate model output y as a function g() of inputs x so that y = g(x). It is trained on an ensemble of model runs at set of inputs called the design matrix, denoted X, giving a sample of model output. The configuration of the design matrix is a latin hypercube (MacKay 1979), as used in e.g. Gregoire et al., 2010, Williams et al., 2013), with sample input points chosen to fill input parameter space efficiently and therefore sample relationships between input parameters effectively.

## 3.1 An augmented emulator

Our strategy is to augment the design matrix of input parameters X with corresponding atmospheric climate model output that might have an impact on the modelled land surface. We build an emulator that models the effects of both input parameters and climate on forest fraction. We then use the augmented emulator to bias correct each forest in turn.

We will use the emulator to describe the relationship between land surface parameters, atmospheric variables that summarise the action of hidden atmospheric parameters, and the broadleaf forest fraction. The relationships between these variables are summarised in figure 4.

We have a number of forests for each ensemble member, differing in driving influences by a different local climate. We use regional mean temperature, $T$, and precipitation, $P$, for each of the forests the Amazon, central Africa and Southeast Asia as additional inputs to augment our original design matrix of land surface parameters, X. Regional extent of each of the broadleaf forests can be found in the supplementary material.

These new inputs are outputs of the model when run at the original inputs $X$, and are influenced by the 10 atmospheric parameters perturbed in a previous ensemble, summarised in the "beta" parameter. We cannot control them directly and thus ensure that they lie in a latin hypercube configuration. We do however hope that they represent a wide spread of model behaviour, given wide perturbations of the input parameters.

With $n = 100$ ensemble members, we form each $n \times 1$ vector of temperature and precipitation and form an $n \times 2$ matrix of climate variables for the Amazon $C_{AZ} = [T_{AZ}P_{AZ}]$, Central Africa $C = [T_{AF}P_{AF}]$ and Southeast Asia $C_{AS} = [T_{AS}P_{AS}]$.

We use these to augment the original $n \times p$ input matrix $X$, creating a unique input location for each forest. We then stack these augmented input matrices together to form a single input matrix X'.

$$X' = \begin{bmatrix} X & C_{AZ} \\ X & C_{AF} \\ X & C_{AS} \end{bmatrix} \tag{1}$$

From an initial ensemble design matrix with $n = 100$ members and $p = 7$ inputs, we now have a design with $n = 300$ members and $p = 9$ inputs. Each member with a replicated set of initial input parameters (e.g members $[1, 101, 201]$), differ only in the $T$ and $P$ values. Figure 5 shows the composition of the resulting input matrix and output vector.

Where in M16, we built an independent emulator for each output (i.e. regional forest fraction), we now build a single emulator for all forest fractions simultaneously given input parameters, temperature and precipitation. The output vector y for the tropical forests has gone from being 3 sets of 100 values, to a single vector $[y_1, \ldots, y_n]$ of length 300.

**Table 1.** Mean absolute error (MAE) rounded to the first significant figure for the regular emulator, using just the seven land surface inputs, and the augmented emulator, including temperature and precipitation.

| Forest | Regular emulator MAE | Augmented emulator MAE |
|---|---|---|
| Amazon | 0.05 | 0.03 |
| Southeast Asia | 0.06 | 0.03 |
| Central Africa | 0.06 | 0.03 |
| All | 0.06 | 0.03 |

We model forest fraction $[y_1, \ldots, y_n]$ as a function of $X'$ using the Gaussian process emulator of package DiceKriging (**?**) in the R statistical language and environment for statistical computing. Details of the emulator can be found in the supplementary material.

## 3.2 Verifying the augmented emulator

To verify that the augmented emulator adequately reproduces the simulator behaviour, we use a leave-one-out metric. For this metric, we sequentially remove one simulator run from the ensemble, train the emulator on the remaining ensemble members and predict the held-out run. We present the predicted members and the calculated uncertainty plotted against the actual ensemble values in figure **?**.

We see no reason to doubt that the emulator provides a good prediction and accurate uncertainty estimates for prediction at inputs points not yet run. We use the mean of the absolute value of the difference between the emulator prediction and corresponding held-out value to calculate the Mean Absolute Error of cross-validation prediction (MAE). Prediction error and uncertainty estimates remain approximately stationary across all tropical forests and values of forest fraction. The mean absolute error of prediction using this emulator is a little under 0.03, or 3% of the maximum possible value of the ensemble.

When compared against the regular emulator using just the land surface inputs, the augmented emulator performs well. The augmented emulator has a mean absolute error of prediction of 0.03 or 3% of the maximum possible value of the ensemble. The regular emulator built individually for each of the forests has a mean absolute error value of 0.058 - nearly double that of the augmented emulator. A breakdown of the mean absolute error of the emulator on a per-forest basis can be seen in table 1

We test the reliability of uncertainty estimates of the emulator by checking that the estimated probability distributions for held-out ensemble members match the true error distributions in the leave-one-out exercise. We create a rank histogram (see e.g. Hamill (2001) **?**) for predictions, sampling 1000 times from each Gaussian prediction distribution, and plotting the rank of the actual prediction in that distribution. The distribution of these ranks overall predictions should be uniform if the uncertainty estimates are reliable. Consistent overestimation of uncertainty will produce a peaked histogram, while systematic underestimation of uncertainty will produce a u-shaped histogram. The rank histogram produced by this set of predictions (figure 7) is close to a uniform distribution, indicating reliable predictions.

## 4 Analyses

### 4.1 Sensitivity analysis

The emulator allows us to measure the sensitivity of forest fraction to the land surface input parameters simultaneously with climate variables temperature and precipitation. We measure the one-at-a-time sensitivity to parameters and climate variables, using the emulator to predict changes in forest fraction as each variable is changed from the lowest to highest setting in turn, with all other parameters at the default settings or observed values. We present the results in figure 8. Parameters NL0 and V_CRIT_ALPHA and climate variables temperature and precipitation exert strong influences of similar magnitudes on forest fraction. Shaded regions represent the uncertainty of the sensitivity to each parameter, due to estimated emulator uncertainty of $\pm 2$ standard deviations. This sensitivity measure does not include the extra uncertainty due to the fact that the relationships will change depending on the position of the other parameters. We do however get to see a measure of how temperature and precipitation affect the marginal response of the other parameters, as the observed climates of each forest are different. We clearly see that the response of the forest fraction to e.g. NL0 depends on climate - the forests fraction response is a noticeably different shape when varied under the mean climate of the South East Asian region.

A quantitative measure of sensitivity of the model output to parameters that does take into account interactions with other parameters is found using the FAST99 algorithm of Saltelli et al. (1999), summarised in figure 9. Precipitation and Temperature are the second and third most important parameters, more important than NL0, and only slightly less important than V_CRIT_ALPHA. Interaction terms contribute a small but non-negligible part to the sensitivity.

### 4.2 The joint impacts of temperature and precipitation on forest fraction

What impact do temperature and precipitation have on forest fraction together? We use the emulator from section **??** and predict the simulator output across the entire range of simulated temperature and precipitation, while holding the other inputs at their default values. The marginal impacts of temperature and precipitation on forest fraction are clear in figure **??**. Ensemble member temperature, precipitation and forest fraction, taken from figure 3 are overplotted for comparison. Temperature and precipitation values are normalised to the range of the ensemble in this plot.

Cooler, wetter climates are predicted to increase forest fraction and drier, warmer climates lead to very low forest fraction values in some simulations. In general, South East Asian and Central African forests are simulated as warmer and wetter than their true-life counterparts. Moving any ensemble member to observed values of temperature and precipitation would not cross many contours of forest fraction value, and so ensemble members are simulated with a roughly accurate forest fraction. In contrast, the Amazon is simulated slightly drier, and considerably warmer than the observed Amazon and many ensemble members consequently have a lower forest fraction than observed. This figure provides strong evidence that a significant fraction of the bias in Amazon forest fraction is caused by a bias in simulated climate.

### 4.3 A climate bias correction approach

With an emulator that models the relationship between input parameters, local climate and the forest fraction, we can predict what would happen to forest fraction in any model simulation if the local climate was correct. In figure 11, for example, the predicted value is the forest fraction at the default set of land surface parameters, with the local temperature and precipitation corrected to the observed values. Central Africa becomes significantly drier, and a little cooler than the centroid of the ensemble. Southeast Asia becomes a little cooler and a little drier. The Amazon forest becomes a little wetter, and significantly cooler. The ensemble has a much larger spread of climates in central Africa than South East Asia or the Amazon.

The bias correction reduces the difference between the modelled and observed Amazon forest fraction markedly, from -0.28 to -0.08. It makes the modelled forest in central Africa worse (-0.11 from -0.03), and slightly improves the SE Asian forest fraction (0.07 from 0.1). Overall, bias correcting the climate takes the mean absolute error at the default parameters from 0.14 to 0.09.

### 4.4 History matching to learn about model discrepancy

In this section we use history matching to learn about parts of input parameter space that are consistent with observations, and to find the causes of discrepancy in the model.

History matching measures the statistical distance between an observation of a real-world process, and the emulated output of the climate model at any input setting. An input where the output is deemed too far from the observation is ruled "implausible", and removed from consideration. Remaining inputs are conditionally accepted as "Not Ruled Out Yet" (NROY), recognising that further information about the model or observations might yet rule them as implausible.

Observations of the system are denoted $z$, and we assume that they are made with uncorrelated and independent errors $\epsilon$ such that

$$z = y + \epsilon \tag{2}$$

Assuming a "best" set of inputs $x^*$ where the the model discrepancy $\delta$, or difference between climate model output $y$ and $z$ is minimised, we relate observations to inputs with

$$z = g(x^*) + \delta + \epsilon \tag{3}$$

We calculate measure of Implausibility I, and reject any input as implausible where I >3 after Pukelsheim's three-sigma rule; that is, for any unimodal distribution, 95 % of the probability mass will be contained within 3 standard deviations of the mean (Pukelsheim, 1994).

We calculate

$$I^2 = |z - E[g(x)]|^2 / Var[g(x)] + Var[\delta]] + Var[\epsilon] \tag{4}$$

**Table 2.** .

| Forest | Error | Implausibility | Error (bias corrected) | Implausibility (bias corrected) |
|---|---|---|---|---|
| Amazon | 0.316 | 6.94 | -0.079 | 1.31 |
| Southeast Asia | -0.096 | 1 .61 | 0.072 | 1.76 |
| Central Africa | -0.04 | 0.768 | -0.11 | 1.5 |

Which recognises that the distance between the best estimate of the emulator and the observations must be normalised by uncertainty in the emulator $g(x)$, in the observational error $\epsilon$, and in the estimate of model discrepancy $\delta$.

We study the region of input space that is "Not Ruled Out Yet" (NROY) by comparison of the model output to the observations of forest fraction. In the previous section we see that the overall difference between the simulated and observed forest fraction is reduced if the output is bias corrected. In this section, we study how that bias correction affects the NROY space.

In M16, the default input parameters were ruled out as implausible for the Amazon region forest fraction. For the sake of illustration, we assume very low uncertainties: zero observational uncertainty and a model discrepancy term with a zero mean and an uncertainty ($\pm$ 1 sd) of just 0.01. We note that under these conditions the default parameters would be ruled out in the standard emulator. However, if we bias correct the model output using the observed temperature and precipitation, we find that the implausibility measure I for the forest fraction in the Amazon at the standard input parameters reduces from nearly 7 to 1.3 - comfortably under the often-used threshold of 3 for rejection of an input. The implausibility of the SE Asian and Central African forest fraction at the default parameter settings rises with bias correction (see table 2), but neither comparison comes close to ruling out the default parameters. We can confidently say that bias correction using the emulator means that observations no longer rule out the default parameters, even with the assumption of a very small model discrepancy.

Another result of bias correction is that it increases the "harmonisation" of the input spaces - that is, the volume of the input space that is "shared", or Not Ruled Out Yet by any of the comparisons of the simulated forest fractions with data. In M16, we argued that the regions of input parameter space where the model output best matched the observations had a large shared volume for The central African, Southeast Asian and North American forests. In contrast, the "best" input parameters for the Amazon showed very little overlap with these other forests. This pointed to a systematic difference between the Amazon and the other forests that might be a climate bias, or a fundamental discrepancy in the land surface component of the model. Here, we show that the climate-bias forest in the Amazon would share a much larger proportion of its NROY space with the other forests. Indeed, the default parameters are now part of this "shared" space, and there is formally no need to invoke an unexplained model discrepancy in order to accept them for all the tropical forests. We offer a cartoon of the situation in figure 12

We find that when we bias correct all the spaces, the proportion of "shared" NROY input space increases from 2.6% to 31% - an order of magnitude increase. This is driven chiefly by the harmonisation of the NROY space of the Amazon to the other two forests. We see that before bias correction, the South East Asian and African forests share nearly ? (74%) of their combined

**Table 3.** .

| Non bias-corrected | Amazon | Southeast Asia | Central Africa |
|---|---|---|---|
| Amazon | 1 | | |
| Southeast Asia | 0.034 | 1 | |
| Central Africa | 0.075 | 0.741 | 1 |

**Table 4.** .

| Bias-corrected | Amazon | Southeast Asia | Central Africa |
|---|---|---|---|
| Amazon | 1 | | |
| Southeast Asia | 0.329 | 1 | |
| Central Africa | 0.915 | 0.337 | 1 |

NROY space. This drops to 33% when bias corrected, but with the advantage that the Amazon and Central Africa now share over 90% (91.5%) of their combined NROY space.

When compared to the initial input parameter space covered by the ensemble, the shared NROY space of the non-bias-corrected forests represents 1.9%, rising to 28% on bias correction.

We visualise two dimensional projections of the NROY input parameter space shared by all three forests before and after bias correction in figure 13 and 14 . The two dimensional projections of high density regions of NROY points are dramatically shifted and expanded in the bias corrected input space, and the default parameters now lie in a high density region.

### 4.5    Doing better than the default parameters

We can use the emulator to find potential locations in parameter space where the difference between the modelled and observed forest fractions would be smaller than at the default parameters. Figure **??** shows the density of parameter settings in each 2-dimensional projection of the input space, where the emulator estimates the model performs better than at the default parameters, once bias correction has been applied. That is, the absolute difference between each estimated forest fraction and the observed values is smaller than the absolute difference of the mean estimate at the default parameters. Out of 100000 samples from the uniform hypercube defined by the range of the experiment design, only 2451, or around 2.5% match this criterion and are plotted. This diagram might help guide further runs in the ensemble, choosing high density regions to run new ensemble members. The convergence of NL0 and V_CRIT_ALPHA seems particularly focussed, and suggests that a lower value of V_CRIT_ALPHA might be a way to reduce error in the forest fraction. There is another, although less densely populated region of high NL0 and V_CRIT_ALPHA that might fulfil the criteria of lower estimates of error for each forest. These regions would be good targets for supplementary runs of the climate model, and for particularly careful emulator checking.

**Table 5.** .

| | Intersection / Union (%) | Intersection / Initial (%) |
|---|---|---|
| Non bias-corrected | 2.6 | 1.9 |
| Bias-corrected | 31 | 28.3 |

## 4.6 Allowable climate at default parameters

We use history matching to find the set of regional mean climates that are most consistent with the observations for each tropical forest. Using the same values for model discrepancy (0), its associated uncertainty (1sd = 0.01), and observational uncertainty (0), we find the set of not-ruled-out-yet temperature and precipitation values when the remaining input parameters are held at their default values. Figure 16 shows the density of not-ruled-out-yet points in the climate space for each of the observed forest fractions. We see that the Amazon and Central African forests might be well simulated in the model in a very wide range of cooler and wetter climates, with only the "hot, dry" corner showing zero density of potential inputs that produce similar forest fraction to observations. The Southeast Asian forest fraction is matched by a swathe of inputs running diagonally through the centre of input space. Neither the hot-dry or cool-wet corners of input space produce forests that match the observations, though the warm-wet and cool-dry corners do.

## 5 Discussion

### 5.1 Simulating the Amazon

What does our study mean for the simulation of the Amazon rainforest? We have shown that the simulation of the broadleaf tropical forest in FAMOUS is almost as sensitive to temperature and precipitation as to any land surface parameter perturbation in the ensemble. It should be remembered however that the calculated sensitivities are dependent on the chosen limits of the parameter perturbations themselves. The precise order and size of sensitivities might change given updated parameter ranges. There is little doubt that the climate variables are a strong influence on broadleaf forest fraction, however.

This version of FAMOUS with the default parameter settings would successfully simulate the Amazon rainforest to within tolerable limits if regional climate biases were substantially reduced. As such, there is no need to invoke a missing process or processes in the land surface in order to explain the discrepancy. Further, there is no need to doubt that the default parameters are implausible and we have strengthened the case made by M16 that the low Amazon forest fraction is a result of poorly chosen parameters. There are also indications from the emulator that a small region of parameter space exists where there is even smaller overall error in the simulation, offering a target for exploration using further runs of the model.

There is a feedback from the land surface to the atmosphere implicitly included in the emulated relationship. We cannot control this feedback directly with the emulator, and so work out the impact of this feedback on the forest fraction as it is

present in the training data. This feedback would have to be taken into account if we were to simulate the correct climate independently of the land surface.

It is possible that were we to include a process seen to be missing from the Amazon (such as deeper rooting of trees allowing them to thrive in drier climates), our map of NROY input space would alter again. Given that there is a measure of uncertainty in observations and the emulator, as well as the possibility of further compensating errors, we cannot rule out a model discrepancy such as a deep rooting process. The fact that the other forests do slightly less well when their climates are bias corrected points to a potential missing process in the model, compensated for by parameter perturbations. However, the impact of this missing process is likely much smaller than we might have estimated had we not taken the bias correction of the forest into account.

## 5.2 Uses for an augmented emulator

By building an emulator that includes temperature and precipitation - traditionally used as climate model outputs - we are able to separate the tuning of one component of the model (the atmosphere) from another (the land surface). Perturbations to the atmospheric parameters, tested in a previous ensemble but not available to us except through an indicator parameter are summarised as inputs through the climate of the model.

We have used the augmented emulator as a translational layer between components of the model. The augmented emulator allows us to ask "what would it mean for our choice of input parameters if the climate of the model in the Amazon region were correct?"This means that we will have less chance of ruling out parts of parameter space that are in fact perfectly good, or keeping those parts that are in fact implausible.

This sort of translational layer might be built as part of a model development process, making it computationally cheaper and faster. Traditionally, the components of computationally expensive flagship climate models are built and tuned in isolation before being coupled together. The act of coupling model components can reveal model discrepancies or inadequacies, and models sometimes have to be re-tuned to work in a coupled mode. There is a danger that this retuning leads to a model that reproduces historical data fairly well, but that makes errors in fundamental processes and therefore is less able to predict or extrapolate - for example, a climate model when projecting future changes under unprecedented greenhouse gas concentrations. Given the time and resources needed to run such complex models, these errors might persist much longer than necessary, and have profound consequences for climate policy. The translational layer allows parameter choices to be made for a model when run in coupled mode, even when there was a significant bias in one of the components that would affect the other components. Using the augmented emulator could could eliminate some of the steps in the tuning process, help the model developer identify potential sources of bias, and to quickly and cheaply calculate the impacts of fixing them. In doing so it would aid model developers in identifying priorities for and allocating effort in future model development.

Our work here shows this process as an example. We have identified the importance of precipitation and temperature to the correct simulation of the Amazon forest, and flag them as priorities in future climate model development. We have identified regions of the space of these climate variables where the Amazon forest might thrive, and related that back to regions of land surface parameter space that might be targeted in future runs of the model.

An augmented emulator could also be used directly to estimate the impacts and related uncertainty of climate change on forest fraction in the model, even in the presence of a significant model bias. After estimating the relationship between the uncertain parameters, climate, and the forest fraction, we could calculate the forest fraction at any climate, including those that

5 might be found in the future. This would require a new ensemble of climate model runs, projecting the future forests under a number of (e.g.) atmospheric CO2 concentrations and parameter combinations. Any extra inputs that might impact the forest, such as atmospheric CO2 concentration, would need to be added to the training set of the emulator. It would be necessary that the training data included any climates that might be seen under the climate change scenario to be studied, as the emulator has much larger uncertainties if asked to extrapolate beyond the limits of the training data. The trajectory of vegetation states

10 through time would also be an important element of the ensemble, as the vegetation state is path dependent. However there would be great potential to save a large number of runs, as not every parameter perturbation would have to be run with every projection scenario.

Such a set of runs would serve as a framework upon which a great many post hoc analyses could be done with the emulator. Once the set of runs was complete, they would effectively serve as the definitive version of the model - any new information

15 that needed to be extracted from the model could in theory be found using the emulator. Not only might we be able to identify and correct important climate biases and their impact on the forest, but also update our estimates of forest change as we learned more about the uncertainty ranges of the uncertain parameters and forcing trajectory.

In theory the augmented emulator could be used to bias correct, for example, every gridbox in a field for a particular variable. The computational resources needed to fit a Gaussian process emulator when the number of outputs estimated simultaneously

20 becomes even moderately large limits the use of our technique. The design input parameter matrix used for training the emulator grows to $n \times d$ rows, where n is the number of ensemble members in the original training set, and d is the number of separate output instances to be considered. In our example, d is 3, and so we only have $100 \times 3 = 300$ in the new training set. Given an initial ensemble of a few hundred, this could easily result in a training set with hundreds of thousands or even millions of rows. Gaussian process emulators are currently limited to using training data with perhaps a few hundred rows as current software

25 packages must invert an $n \times n$ matrix, a potentially very computationally expensive process (see e.g. Hensman et al. (2013) for examples). At the time of writing this limitation would preclude using our specific technique for correcting biases on a per-gridbox basis. To make use of the translational layer for large data sets we would need new Gaussian process technology, or specific strategies to deal with large data sets. These strategies might involve kernel based methods, keeping the scope of training data local to limit the size of any inverted matrices. Alternatively, they might involve building emulators using only a

30 strategically sampled selection of the outputs. Recent advances in using Gaussian processes for larger data sets can be found in Hensman et al. (2015, 2013), Wilson et al. (2015a) 2015b.

Given that we overcome such technical barriers, we see no reason that such a layer not be built that is used to (for example) correct the climate seen by individual land surface grid boxes, rather than (as here) individual aggregated forests. The process of rejecting poor parameter sets might be aided by having a comparison against each gridbox in an entire global observed surface,

35 rather than aggregated forests. Alternatively, we might allow parameters to vary on a gridbox-by-gridbox basis, effectively forming a map of Not-Ruled-Out-Yet parameters.

There are also potential computational efficiencies in our approach of decoupling the tuning of two components (here the atmosphere and the land surface) in the model. A good rule of thumb is that a design matrix for building an emulator should have $O(10 \times p)$ training points, where p is the number of input parameters, in order to adequately sample parameter space to the extent it is possible to build a good emulator. With approximately 10 atmospheric and 7 land surface parameters, we would need O(170) runs. Here, we have summarised those 10 parameters as two outputs that have a material impact on the aspect of the land surface that we are interested in. Adding these two to the 7 inputs, we need $O(10 \times (2 + 7) = 90)$ runs, well covered by our available ensemble of 100 runs.

## 6    Conclusions

We demonstrate that we can correct the simulation of the Amazon rainforest in the climate model FAMOUS by correcting the regional bias in the climate of the model with a Gaussian process emulator. We therefore find it unnecessary to invoke a model discrepancy or inadequacy, such as a lack of deep rooting in the Amazon in the model, to explain the anomalously low forest fraction in an ensemble of forests simulations.

We present a method of augmenting a Gaussian process emulator by using climate model outputs as inputs to the emulator. We use average regional temperature and precipitation as inputs, alongside a number of land surface parameters, to predict average forest fraction in the tropical forests of the Amazon, Southeast Asia and central Africa. We assume that the differences in these parameters account for the regional differences between the forests, and use data from all three tropical forest regions to build a single emulator.

We find that the augmented emulator improves accuracy in leave-one-out test of prediction, reducing the mean absolute error of prediction by almost half, from nearly 6% of forest fraction to just under 3%. This allays any fears that the emulator is inadequate to perform a useful analysis, or produces a measurable bias in predictions, once augmented with temperature and precipitation as inputs.

In two types of sensitivity analyses, temperature and precipitation are important inputs, ranking 2 and 3 after V_CRIT_ALPHA (rank 1) and ahead of NL0 (rank 4).

We use the augmented emulator to bias correct the climate of the climate model to modern observations. Once bias corrected, the simulated forest fraction in the Amazon is much closer to the observed value in the real world. The other forests also change slightly, with central Africa moving further from the observations, and Southeast Asia moving slightly closer. We find that the differences in the accuracy of simulation of the Amazon forest fraction and the other forests can be explained by the error in climate in the Amazon.

There is no requirement to invoke a land surface model discrepancy in order to explain the difference between the Amazon and the other forests. After bias correction, the default parameters are classified as "Not Ruled Out Yet" in a history matching exercise, that is they are conditionally accepted as being able to produce simulations of all three forests that are statistically sufficiently close to the values observed in the real world.

17

Bias correction "harmonises" the proportion of joint NROY space that is shared by the three forests. This proportion rises from 2.6% to 31% on bias correction.

Taken together these finding strengthen the conclusion of McNeall et al. 2016 that the default parameters should not be ruled out as implausible by the failure of FAMOUS to simulate the Amazon.

We find a small proportion (around 2.5%) of input parameter space where we estimate that the climate model might simulate the forests better than at the default parameters. This space would be a good target for further runs of the simulator.

Finally, we offer a technique of using an emulator augmented with input variables that are traditionally used as outputs, to aide the tuning of a coupled model perturbed parameter ensemble by separating the tuning of the individual components. This has the potential to (1) reduce the computational expense by reducing the number of model runs needed during the model tuning and development process and (2) help model developers prioritise areas of the model that would most benefit from development. The technique could also be applied to efficiently estimate the impacts of climate change on the land surface, even where there are substantial biases in the current climate of the model.

*Code availability.* TEXT

*Data availability.* TEXT

*Code and data availability.* TEXT

*Sample availability.* TEXT

*Video supplement.* TEXT

## Appendix A

## A1

*Author contributions.* TEXT

# References

Andrianakis, I., Vernon, I. R., McCreesh, N., McKinley, T. J., Oakley, J. E., Nsubuga, R. N., Goldstein, M., and White, R. G.: Bayesian history matching of complex infectious disease models using emulation: a tutorial and a case study on HIV in Uganda, PLoS computational biology, 11, e1003 968, 2015.

Arendt, P. D., Apley, D. W., and Chen, W.: Quantification of model uncertainty: Calibration, model discrepancy, and identifiability, Journal of Mechanical Design, 134, 100 908, 2012a.

Arendt, P. D., Apley, D. W., Chen, W., Lamb, D., and Gorsich, D.: Improving identifiability in model calibration using multiple responses, Journal of Mechanical Design, 134, 100 909, 2012b.

Brynjarsdóttir, J. and O'Hagan, A.: Learning about physical parameters: the importance of model discrepancy, Inverse Problems, 30, 114 007, http://stacks.iop.org/0266-5611/30/i=11/a=114007, 2014.

Craig, P., Goldstein, M., Seheult, A., and Smith, J.: Bayes linear strategies for history matching of hydrocarbon reservoirs, in: Bayesian Statistics 5, edited by Bernardo, J., Berger, J., Dawid, A., and Smith, A., pp. 69–95, Clarendon Press, Oxford, UK, 1996.

Golaz, J.-C., Horowitz, L. W., and Levy, H.: Cloud tuning in a coupled climate model: Impact on 20th century warming, Geophysical Research Letters, 40, 2246–2251, 2013.

Hourdin, F., Mauritsen, T., Gettelman, A., Golaz, J.-C., Balaji, V., Duan, Q., Folini, D., Ji, D., Klocke, D., Qian, Y., et al.: The art and science of climate model tuning, Bulletin of the American Meteorological Society, 98, 589–602, 2017.

Kennedy, M. and O'Hagan, A.: Bayesian calibration of computer models, Journal of the Royal Statistical Society: Series B (Statistical Methodology), 63, 425–464, 2001.

McNeall, D., Williams, J., Booth, B., Betts, R., Challenor, P., Wiltshire, A., and Sexton, D.: The impact of structural error on parameter constraint in a climate model., Earth System Dynamics, 7, 2016.

McNeall, D. J., Challenor, P. G., Gattiker, J. R., and Stone, E. J.: The potential of an observational data set for calibration of a computationally expensive computer model, Geoscientific Model Development, 6, 1715–1728, https://doi.org/10.5194/gmd-6-1715-2013, http://www.geosci-model-dev.net/6/1715/2013/, 2013.

Schmidt, G. A., Bader, D., Donner, L. J., Elsaesser, G. S., Golaz, J.-C., Hannay, C., Molod, A., Neale, R. B., and Saha, S.: Practice and philosophy of climate model tuning across six US modeling centers, Geoscientific Model Development, 10, 3207–3223, 2017.

Vernon, I., Goldstein, M., and Bower, R.: Galaxy formation: a Bayesian uncertainty analysis, Bayesian Analysis, 5, 619–669, 2010.

Walters, D., Baran, A., Boutle, I., Brooks, M., Earnshaw, P., Edwards, J., Furtado, K., Hill, P., Lock, A., Manners, J., Morcrette, C., Mulcahy, J., Sanchez, C., Smith, C., Stratton, R., Tennant, W., Tomassini, L., Van Weverberg, K., Vosper, S., Willett, M., Browse, J., Bushell, A., Dalvi, M., Essery, R., Gedney, N., Hardiman, S., Johnson, B., Johnson, C., Jones, A., Mann, G., Milton, S., Rumbold, H., Sellar, A., Ujiie, M., Whitall, M., Williams, K., and Zerroukat, M.: The Met Office Unified Model Global Atmosphere 7.0/7.1 and JULES Global Land 7.0 configurations, Geoscientific Model Development Discussions, 2017, 1–78, https://doi.org/10.5194/gmd-2017-291, https://www.geosci-model-dev-discuss.net/gmd-2017-291/, 2017.

Williamson, D., Blaker, A. T., Hampton, C., and Salter, J.: Identifying and removing structural biases in climate models with history matching, Climate dynamics, 45, 1299–1324, 2015.

Zhao, M., Golaz, J.-C., Held, I., Guo, H., Balaji, V., Benson, R., Chen, J.-H., Chen, X., Donner, L., Dunne, J., et al.: The GFDL global atmosphere and land model AM4. 0/LM4. 0: 2. Model description, sensitivity studies, and tuning strategies, Journal of Advances in Modeling Earth Systems, 10, 735–769, 2018.
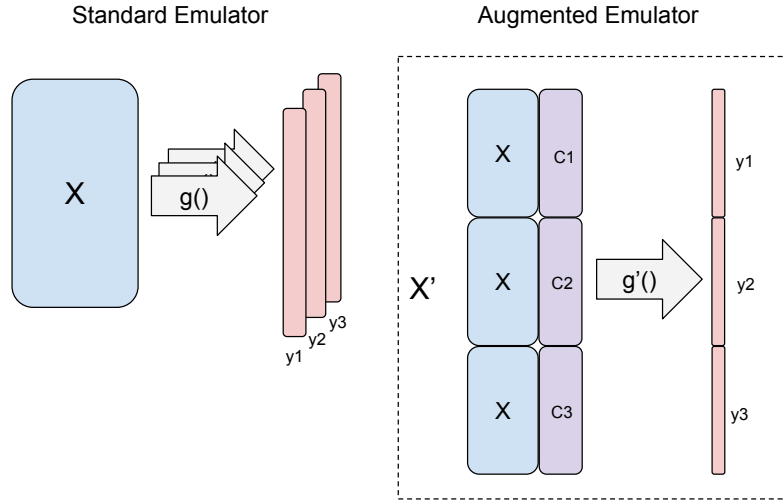
**Figure 2.** Observations of broadleaf forest fraction on their native grid (top), and regridded to the FAMOUS grid (bottom)

**Figure 3.** Smaller symbols represent broadleaf forest fraction in the FAMOUS ensemble against regional mean temperature and precipitation. Ensemble member forest fraction in the Amazon is represented by the colour of the circles, Central Africa by triangles and SE Asia by squares. Larger symbols represent observed climate and forest fraction.
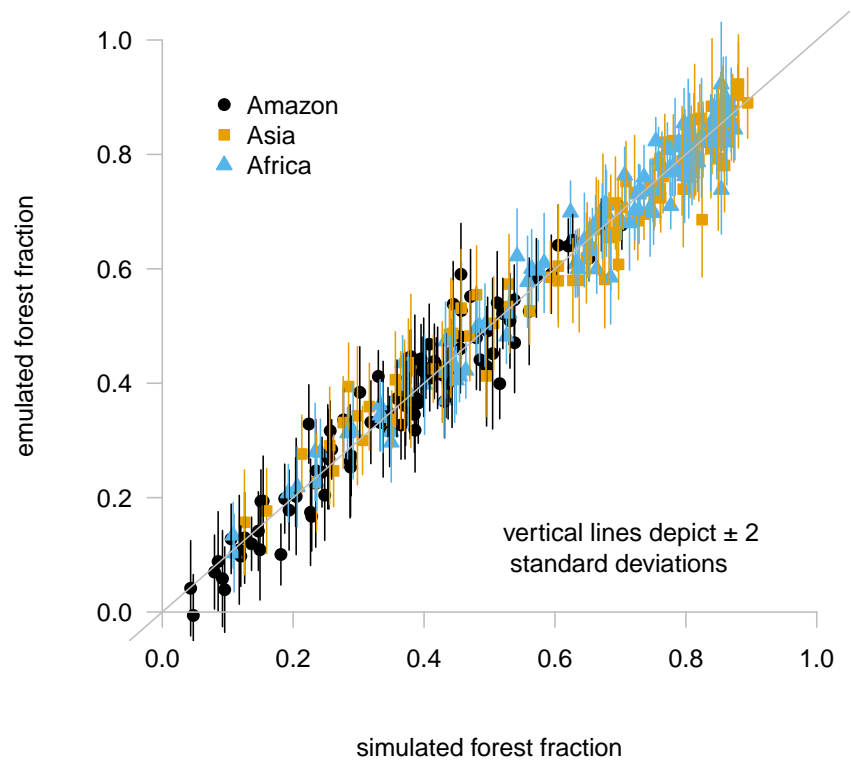
**Figure 4.** A graph showing the assumed relationship between input parameters, climate variables and forest fraction. An arrow indicates influence in the direction of the arrow. Processes that are directly emulated are shown with a solid arrow, while the processes shown by a dotted arrow are not directly emulated.

**Figure 5.** In a standard emulator setup (left), training data consists of an input matrix $X$ and corresponding simulator output $y$. A new emulator $g_1, \ldots, () g_n()$ is trained for each output $y_1, \ldots, y_n$ of interest. In the augmented emulator, output from the simulator $C_1, \ldots, C_3$ augments the design matrix, with the initial inputs $X$ repeated.

**Figure 6.** Leave-one-out cross validation plot, with the true value of the simulator output on the x-axis, and predicted output on the y-axis. Vertical lines indicate $\pm 2$ standard deviations.

**Figure 7.** Rank histogram of leave-one-out predictions. For each prediction of a held-out ensemble member, we sample 1000 points from the Gaussian prediction distribution, and then record where the true held-out ensemble member ranks in that distribution. We plot a histogram of the ranks for all 300 ensemble members. A uniform distribution of ranks indicates that uncertainty estimates of the emulator are well calibrated.

**Figure 8.** One-at-a-time sensitivity of forest fraction variation of each parameter and climate variable in turn across the entire ensemble range. All other parameters or variables are held at their default values while each parameter is varied. Solid lines represent the emulator mean and shaded areas represent $\pm$ 2 standard deviations of emulator uncertainty.
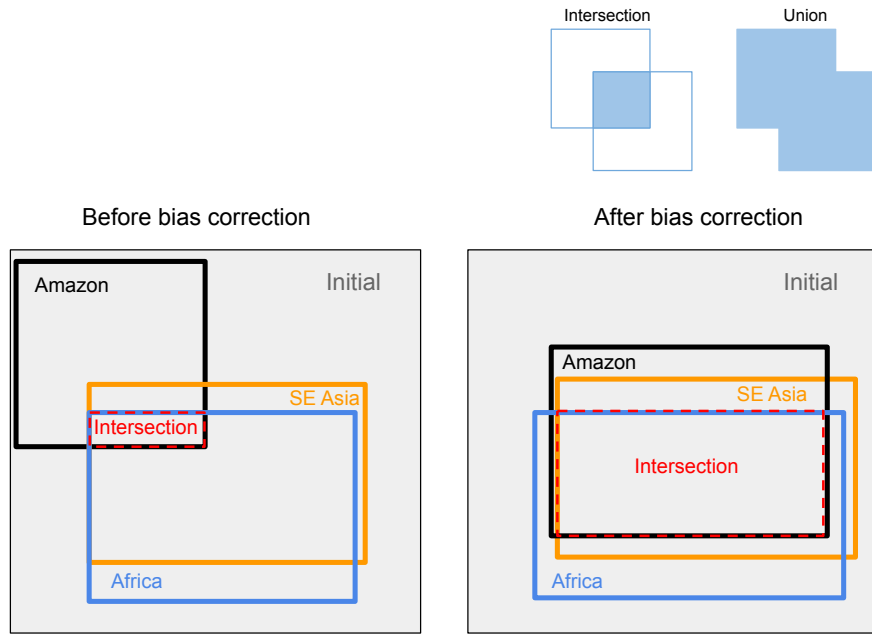
**Figure 9.** Sensitivity of forest fraction to model parameters and climate parameters, found using the FAST99 algorithm of Saltelli et al 1999.

**Figure 10.** The impact of climate on forest fraction. Background plot colour indicates the mean emulated forest fraction when all land surface inputs are held at their default values. Temperature and precipitation in the ensemble are marked with symbols, with the fill colour representing forest fraction. Larger symbols represent the values observed in the real world.
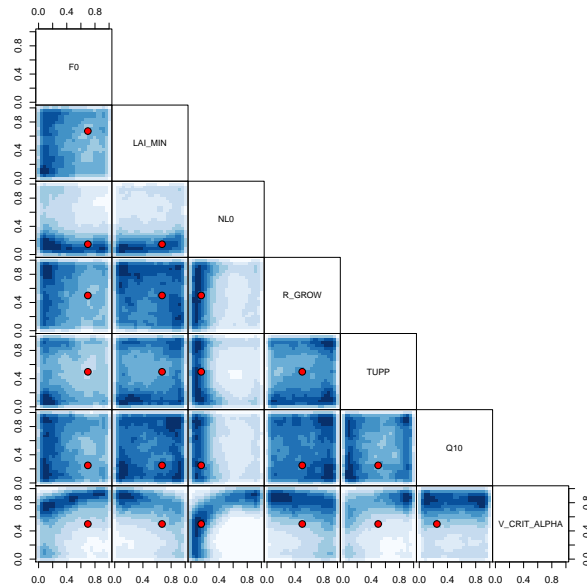
**Figure 11.** Observed and emulated Forest fraction in each tropical forest. For the emulated forest fraction at default and bias corrected parameters, emulator uncertainty of $\pm\,2sd$ is represented by horizontal bars.

**Figure 12.** A cartoon depicting the input space that is "not ruled out yet (NROY)" when the climate simulator output is compared to observations of the forest fraction in the Amazon, Africa, and South East Asia before (left) and after (right) bias correction. We measure the "shared" space (the intersection of NROY spaces for each forest) as a fraction of the union (the total space covered by all three forests) of the NROY spaces. The "initial" space represents the total parameter space covered by the ensemble.
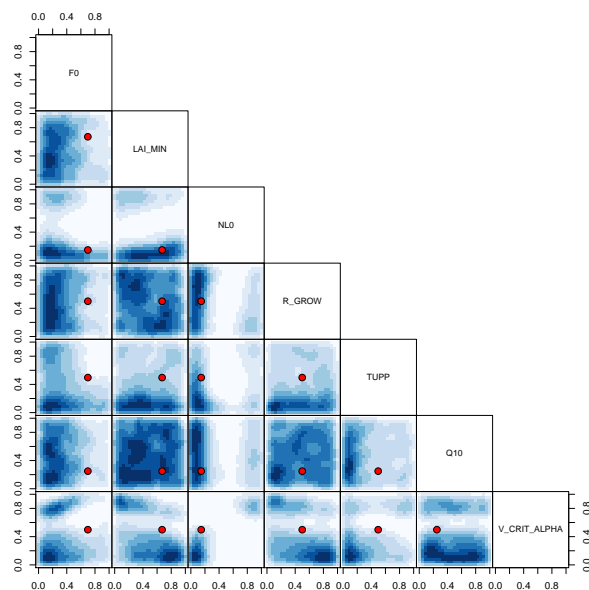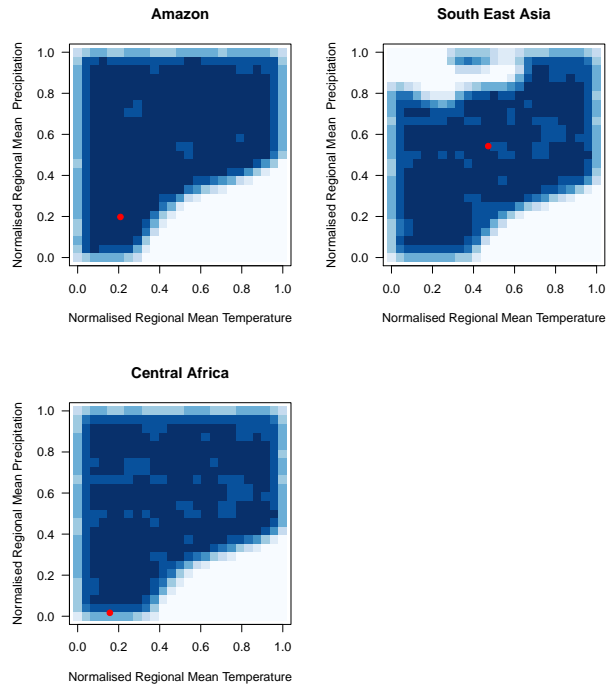
**Figure 13.**



**Figure 14.** Transformation of the NROY land surface input space shared by all three forests. Blue shading denotes the density of NROY input candidates, projected into the two dimensional space indicated by the labels. The default parameter settings are marked as red points. The default parameters are in a low density region in the non-biased-corrected diagram (top), but in a high density region in the bias-corrected diagram (bottom).

**Figure 15.** Two dimensional projections of the density of inputs where the corresponding bias corrected emulated forests have a smaller error than the bias corrected default parameters. These regions might be good targets for additional runs of the climate model. Default parameters are shown as a red point.

**Figure 16.** Density of not-ruled-out yet emulated temperature and precipitation pairs for each observed tropical forest fraction, when input parameters are held at their default values. Observed climates for each forest are marked in red.