

ABC for climate: dealing with expensive simulators

P. B. Holden^{1,*}, N. R. Edwards¹, J. Hensman², and R. D. Wilkinson³

¹Earth, Environment and Ecosystems, Open University, UK

²Computer Science, University of Sheffield, UK

³School of Mathematics and Statistics, University of Sheffield, UK

**Corresponding author*

November 12, 2015

1 Introduction

One of the primary challenges faced when calibrating a simulator using ABC is overcoming the computational constraints posed by working with limited resource. The requirement to repeatedly simulate from a model can make inference extremely computationally expensive. Consequently, much of the methodological development in ABC has focused on improving computational efficiency, either through the use of more efficient Monte Carlo algorithms, or through the use of statistical methods to ameliorate the effect of using a large tolerance.

The difficulty of dealing with limited computer power is felt more keenly in climate science than in most other disciplines. A major focus of climate research concerns the construction of ever more accurate and comprehensive simulators of the climate system. Since the 1970's, global climate models have evolved from representing only the large-scale circulation of the global atmosphere (e.g., Holloway Jr and Manabe, 1971) to models that incorporate complex dynamic representations of land surface, ocean, sea ice, atmospheric aerosols, ocean biogeochemistry, vegetation, soils and atmospheric chemistry (Flato et al., 2013). Separate Earth system components are coupled through

the exchange of fluxes, which describe the flow of some quantity between them (e.g. energy, moisture, CO_2), and by passing any state variables that are needed to define boundary conditions (e.g. land surface albedo, sea surface temperature). “Intermediate complexity” models (which use simplified model components and lower resolution in return for a more complete description of the Earth system and higher computational efficiency) may also include dynamic representations of other important elements, such as ice sheets, permafrost, ocean sediments and weathering (Flato et al., 2013), but these additional, long-timescale components require orders of magnitude longer simulations to reach equilibrium. Modern climate models are generally, and more accurately, described as “Earth system models” or ESMs.

This evolution in complexity has been accompanied by a 5-fold increase in spatial resolution, allowing the resolution of important finer scale processes. This increased resolution (combined with shorter time-steps that are required for numerical stability at higher spatial resolution) has *alone* led to an $O(1000)$ -fold increase in computational demands since the 1970’s. In general higher resolution allows more direct and more realistic representation of smaller-scale processes, although this does not guarantee better projections, in part because more complex models are more challenging to calibrate. A feature of climate modelling is that multi-decadal climate projections must be used before data are available to validate them, while past data give only approximate clues to the expected behaviour of model discrepancy because expected changes greatly exceed the range of variability in the instrumental period.

It is perhaps inevitable, given the continual striving for more complex models and the highest possible resolution, that state-of-the-art ESMs will always be at the limits of what is practicable with available computing power. The UK Met Office Hadley Centre’s computer comprises eight ‘supernodes’ of IBM Power775 supercomputer servers, which were installed in 2012 at a cost of more than £11 million. The ESMs run at the Hadley Center and at equivalent climate modelling institutions in other countries are extremely computationally expensive, requiring months of such supercomputing to perform a single simulation of order 100 years. Even the intermediate complexity model GENIE-1 (Holden et al., 2013a) used in our case study (Section 4) requires several days (on a single CPU node) to perform each $O(10)$ kyear “spin-up” simulation to reach equilibrium, so that simulation ensembles require implementation on multi-node computing clusters. The simulated climates are large complex datasets which comprise temporally-resolved three-dimensional spatial arrays of up to ~ 100 state variables. These outputs, in particular the outputs of carefully designed model inter-comparison projects, are often analysed in great detail, in a comparable way to how scientists in

other fields analyse the outputs from empirical studies; model projections are the best, and only predictions we have of future climate.

An ESM configuration is determined by the settings of many 100’s of model parameters. These include switches (which determine the precise numerical schemes applied), physical constants that are approximately known but vary spatially in the real world (such as the reflectivity of ice) and parameterisations of “sub-gridscale” processes such as cloud formation, which have “tuned” values that are known to result in reasonable model behaviour. This complexity (many weakly constrained inputs, high dimensional outputs and expensive simulators) has meant that careful statistical calibration (either with Bayesian or frequentist approaches) does not have a long history in climate science. Often different modules of an Earth system simulator are separately “tuned” before being bolted together. For example, the atmospheric component can be tuned independently of the ocean component by prescribing sea-surface temperatures with observational values. The components may all be tuned independently before being coupled, with no guarantee that what was a good tuning in an isolated module will work well in the coupled model. After coupling, a small subset of model parameters are adjusted so that the coupled model is consistent with large-scale observational constraints. It has been shown, perhaps unsurprisingly, that such a tuning process does not produce a unique solution, so that different combinations of parameters can lead to equally plausible model realisations (Mauritsen et al., 2012).

As statistical methodology develops, scientists are beginning to perform more careful parameter estimation in their models. More rigorous parameter estimation methods are often developed with (relatively fast) intermediate complexity models, e.g., Annan et al. (2005), thereby informing application to higher-complexity models, e.g., Marquis et al. (2014). We can view climate simulators as black boxes which map from parameter values $\theta \in \Theta$, to climate states $f(\theta) = \mathcal{C}_{\text{sim}}$. The aim of a Bayesian calibration, is to find the posterior distribution

$$\pi(\theta|\mathcal{C}_{\text{obs}}) \propto \int \pi(\mathcal{C}_{\text{obs}}|\mathcal{C}_{\text{sim}})\pi(\mathcal{C}_{\text{sim}}|\theta)d\mathcal{C}_{\text{sim}}\pi(\theta), \quad (1)$$

where \mathcal{C}_{obs} is a set of observations of the climate system (Kennedy and O’Hagan, 2001; Rougier, 2007). Here, $\pi(\theta)$ is the prior distribution for θ , $\pi(\mathcal{C}_{\text{sim}}|\theta)$ is the simulator likelihood function (which is typically unknown) and $\pi(\mathcal{C}_{\text{obs}}|\mathcal{C}_{\text{sim}})$ is the statistical model relating the simulator to physical climate. This calculation, however, is typically far too ambitious to perform in practice. Computational restrictions generally limit us to an ensemble of N simulator runs $\{\theta^{(i)}, \mathcal{C}_{\text{sim}}^{(i)}\}_{i=1}^N$. Typically, N is small, ruling out most Monte

Carlo based calibration approaches. We are left needing to estimate $\pi(\theta|\mathcal{C}_{\text{obs}})$ as best we can, often by adding further approximation.

A further problem faced by climate scientists is that simulator discrepancy (often called model error) can be considerable (Murphy et al., 2004). And whilst the physical models of climate, $\pi(\mathcal{C}_{\text{sim}}|\theta)$, are well developed, statistical models of the simulator discrepancy relating simulated to observed climate, $\pi(\mathcal{C}_{\text{obs}}|\mathcal{C}_{\text{sim}})$, have only begun to be developed relatively recently (Rougier and Goldstein, 2014). The large simulator discrepancy makes most simulators incapable of reproducing all aspects of the climate record simultaneously and can mean that the simulator parameters are no longer directly comparable to their physical namesakes, making prior specification challenging.

So what is possible? We know that ABC, given infinite computational resource and a perfect simulator, can in theory produce arbitrarily accurate posteriors (i.e., the ABC posterior can be made arbitrarily close to the true posterior). But for many problems, computational resource is often severely constrained and simulator discrepancy can be significant and largely unmodelled. Climate science is interesting for the statistician as it presents extreme cases of both these issues.

A key idea allowing calibration in many of these expensive simulators is the idea of replacing the simulator with an emulator (or meta-model), which is a cheap statistical surrogate used in place of the simulator (Sacks et al., 1989; Santner et al., 2003; O’Hagan, 2006). Emulation techniques are attracting considerable interest in the climate community. They are used, for instance, to approximate probabilistic model outputs (Sansó et al., 2008; Rougier et al., 2009; Harris et al., 2013), for parameter estimation (Sham Bhat et al., 2012; Olson et al., 2012), to facilitate model understanding (Lee et al., 2012; Holden et al., in press) and to provide numerically efficient model surrogates for coupling applications (Castruccio et al., 2014; Holden et al., 2014; Oyebamiji et al., 2015). The application we will describe here is in the ABC design of “plausible” simulation ensembles (Holden et al., 2010; Edwards et al., 2011), using emulation in order to overcome the prohibitive limitations imposed by simulator cost.

2 History-matching and ABC

Climate science presents the double whammy of computationally expensive simulators, and simulator discrepancy that is too large to ignore but which is not well understood or modelled. Both of these issues make a careful Bayesian calibration (as described by Equation 1) difficult. What can be achieved? Our aim is to compare observations of Earth’s climate \mathcal{C}_{obs} , with

simulator predictions $\mathcal{C}_{\text{sim}} = f(\theta)$, in order to learn about the parameter θ . ABC is an approach for obtaining a probabilistic calibration, and seeks to match simulator output to observations, approximating the distribution

$$\pi_{ABC}(\theta|\mathcal{C}_{\text{obs}}) \propto \int \mathbb{I}(\rho(\mathcal{C}_{\text{obs}}, \mathcal{C}_{\text{sim}}) \leq \epsilon) \pi(\mathcal{C}_{\text{sim}}|\theta) d\mathcal{C}_{\text{sim}} \pi(\theta). \quad (2)$$

The acceptance kernel $\mathbb{I}(\rho(\mathcal{C}_{\text{sim}}, \mathcal{C}_{\text{obs}}) \leq \epsilon)$ implicitly implies a uniform distribution for the simulator discrepancy (Wilkinson, 2013), but this is usually viewed as a pragmatic compromise, rather than a modelling decision.

An alternative to a probabilistic calibration, is to do a history match (Williamson et al., 2013), which has been used in studies involving complex computer models, such as oil reservoir modelling (Craig et al., 1997), cosmology (Vernon et al., 2010), epidemiology (Andrianakis et al., 2015), and climate science (Edwards et al., 2011). History matching, like calibration, seeks to identify regions of the input space that give acceptable matches between simulator output, \mathcal{C}_{sim} , and observed data, \mathcal{C}_{obs} . But instead of finding a probability distribution over Θ , we instead seek merely to rule out implausible regions of input space, i.e., those θ that the simulator suggests could not have lead to \mathcal{C}_{obs} , even after having accounted for the simulator discrepancy. Often large parts of the input space give simulated climates that are very different from the observed data, and which can hence be ruled to be physically implausible and removed from further consideration.

We define $\mathcal{P}_{\mathcal{C}}$ to be a set of plausible climate states that represent an acceptable match between simulation and observation. We define \mathcal{P}_{θ} to be the subset of the parameter space that leads to plausible simulated climates, i.e.,

$$\mathcal{P}_{\theta} = \{\theta \in \Theta : f(\theta) \in \mathcal{P}_{\mathcal{C}}\}.$$

Often, the vast majority of the input space gives rise to unacceptable matches to the observed data (sometimes $\mathcal{P}_{\theta} = \emptyset$), and it is these regions that we are trying to rule out as implausible. For example, for an ESM, we might define $\mathcal{P}_{\mathcal{C}}$ to be any simulated climate that has global surface air temperature (SAT) within 2°C of the observed value, the maximum value of Atlantic Meridional Overturning Circulation (AMOC), a measure of the large-scale circulation of the ocean, within 5 Sv (1 Sv = $10^6 m^3 s^{-1}$) of observations, and the global mass of vegetation to be within 200 giga-tonnes carbon (GTC) of observations, though clearly the choice of appropriate metrics and acceptance ranges is highly simulator-dependent. \mathcal{P}_{θ} is then the set of model parameters that would generate plausible climates for the ESM in question.

Note the similarity to ABC here. If the prior distribution for θ is uniform on Θ , i.e., $\pi(\theta) \propto \mathbb{I}_{\theta \in \Theta}$, if $f(\theta)$ is deterministic (as is often, at least approximately, the case in climate science), and if we use $\mathbb{I}_{f(\theta) \in \mathcal{P}_{\mathcal{C}}}$ as the ABC

acceptance kernel, then

$$\pi_{ABC}(\theta|\mathcal{C}) \propto \begin{cases} 1 & \text{if } \theta \in \mathcal{P}_\theta \\ 0 & \text{otherwise.} \end{cases}$$

If we interpret a posterior probability of zero, as the statement that θ is implausible, then history matching and ABC are thus the same. Note also the direct relationship between the discrepancy considerations built into $\mathcal{P}_\mathcal{C}$, and the way ABC performs ‘Monte Carlo’ exact inference for the model that has a discrepancy defined by the acceptance kernel (Wilkinson, 2013).

History matching and ABC have in common that they do not use a detailed model of the discrepancy, but instead characterise it using simple criteria. A philosophical difference between the two approaches perhaps lies in the degree of thought given to the plausible set $\mathcal{P}_\mathcal{C}$. In history matching, the plausibility criteria are often based on measurement error variances and the expected magnitude of the simulator discrepancy (Vernon et al., 2010). Consequently, \mathcal{P}_θ consists of those parameter values θ that have not yet been ruled out as implausible by our knowledge of the simulator and its discrepancy, and the observed data and measurement error. The result is usually not interpreted probabilistically, but only as values that we can not yet rule to be implausible given our current state of knowledge. In contrast, in ABC the choice of metric ρ and tolerance ϵ are usually based pragmatically on the characteristics of the algorithm, rather than on physical aspects of the underlying problem. Often, ϵ is chosen to generate a specified number of acceptances. For example, if the computational budget allows for 10^8 simulator runs, and we want 10^4 accepted values in order to approximate the posterior, we set ϵ to the value that leads to 0.01% of simulations being accepted (for example, Biau et al., 2015, interpret ABC as a nearest neighbour algorithm).

A further difference lies in the choice of information to include in $\mathcal{P}_\mathcal{C}$ (i.e., what summary statistics to use). Climate simulators provide a large variety of outputs, and some of these are better able to reproduce observed climate than others. For example, temperatures are generally better reproduced than precipitation, consequently, it is more common to calibrate to the former than the latter. In contrast, ABC has its roots in genetics, where perhaps the simulator output is less varied, and consequently, more focus is given to the automatic selection of summary statistics, often chosen on the basis of what is most informative for θ (Blum et al., 2013). This approach is unlikely to be suitable in climate science. Some outputs for which the simulator discrepancy is particularly large (precipitation say) may well be very ‘informative’ about θ if we do not allow for discrepancy, but this would only misguide and may lead us to incorrectly rule out large swathes of parameter space as implausible.

Variables which are not well simulated are often included in ESMs, either because they improve the overall simulation through the representation of important feedbacks, or because they are considered important outputs in their own right in spite of higher discrepancy associated with the outputs. Whether a weak calibration constraint on these outputs is appropriate will depend on the details of the discrepancy. Where a known missing process gives a significant contribution to regional error for instance, such as large precipitation errors in monsoon regions as a result of unresolved topographic variation, using a too precise calibration constraint (equivalently too small a model discrepancy) would distort the rest of the solution.

A key question for any simulator is whether given a set of plausibility conditions, the simulator is capable of producing any plausible simulated climates. That is, is \mathcal{P}_θ empty? If \mathcal{P}_θ is empty, it is an indication that we understand less than we thought about the simulator and system. Either there is an error in our implementation of the simulator, or we have underestimated the magnitude of the simulator discrepancy or measurement error. The fact that the result of a history-match can be to find there are no plausible parameter values should not be seen as a negative aspect of the approach, as it forces us to confront the cold reality that something is missing from our understanding of the system. In contrast, likelihood based techniques such as MCMC (and pragmatic ABC applications where ϵ is chosen to guarantee a particular acceptance rate), result in a posterior distribution $\pi(\theta|\mathcal{C}_{\text{obs}})$ regardless of how close the simulated climates are to real climate. It is thus sensible when using these techniques, to carefully check that the calibrated simulator does indeed produce acceptable fits. While it can often be useful to find the distribution $\pi(\theta|\mathcal{C}_{\text{obs}})$ (or an approximation to it) regardless of the simulator quality, note that if discrepancy is ignored, $\pi(\theta|\mathcal{C}_{\text{sim}})$ can often be more constrained, or equivalently $|\mathcal{P}_\theta|$ smaller, than is justified (Brynjarsdóttir and O’Hagan, 2014).

Note that even if a probabilistic calibration is required, a history match can be performed first in order to rule out regions of space which are clearly implausible. This can dramatically reduce the area needed to be explored during the more challenging probabilistic calibration. If using a stochastic simulator, for which θ may never be completely ruled as implausible (as $\pi(\theta|\mathcal{C}_{\text{sim}}) > 0$ for all θ say), this can still be advantageous. We can rule out parameter regions for which the likelihood is considerably smaller than at the MLE with only a small increase in the approximation error (Wilkinson, 2014), again making a subsequent probabilistic calibration easier.

Assuming that \mathcal{P}_θ is not empty, the question then becomes, can we find elements of \mathcal{P}_θ , and better still, can we characterise all of \mathcal{P}_θ ? The complexity of climate science is such that even incomplete specifications of \mathcal{P}_θ are useful,

as discussed in Section 4. This is because interest lies not in \mathcal{P}_θ , but in what it implies about future climate, i.e., in the implied calibrated distribution for other aspects of the climate system

$$\pi(\mathcal{C}_{\text{future}}|\mathcal{C}_{\text{obs}}) = \int \pi(\mathcal{C}_{\text{future}}|\theta)\pi(\theta|\mathcal{C}_{\text{obs}})d\theta.$$

and so even partial descriptions of \mathcal{P}_θ are useful in constraining our beliefs about future climate behaviour. Our aim is thus, given a limited computational budget of N simulator evaluations, can we find \mathcal{P}_θ and the corresponding set of plausible future climates? ABC applications usually use millions of simulator evaluations. What can we do if instead we can only afford 100 or 1000 simulator evaluations? The answer is going to be even more approximate than in ABC, and furthermore, we will necessarily have to make some modelling assumptions if we wish to make progress. The key tool that has arisen for doing this is the emulator, or meta-model.

3 Emulation

If the simulator, $f(\theta)$, is expensive to evaluate, we can instead try to find an approximation, $\tilde{f}(\theta)$, called an emulator or meta-model, which provides a good approximation to $f(\theta)$, but which is computationally cheap (Sacks et al., 1989; O’Hagan, 2006). We can then either use \tilde{f} to answer the question of interest (e.g., calibrating the simulator), or use it to guide the choice of the next parameter value at which to evaluate f .

We start by generating an ensemble of simulator evaluations $\mathcal{D} = \{\theta_i, f(\theta_i)\}_{i=1}^N$, which we then use to build \tilde{f} . Building an emulator is a regression problem, and consequently a myriad of different techniques have been used, including linear regression and its variants, neural networks, and Gaussian processes (GPs, also known as Kriging), with GPs proving the most popular class of model thus far. The functional form of the simulator is not known a priori, and so neither is the best regression model, but a reasonable approximation can usually be found using GPs, as long as the response is a smooth continuous function of θ . For the purposes of calibration, the key properties of any emulator are predictive accuracy, quantification of uncertainty in the predictions, and speed of prediction. In climate science, where the output fields being modelled are often spatio-temporal fields, the regression model is usually combined with a dimension reduction technique, to project the output onto a lower dimensional manifold (Higdon et al., 2008; Holden and Edwards, 2010; Wilkinson, 2010).

3.1 Sequential history matching

For many problems, the plausible set \mathcal{P}_θ may constitute only a small fraction of the prior space Θ . Furthermore, \mathcal{P}_θ may consist of multiple disconnected regions. For Monte Carlo methods, this can make designing an effective sampler difficult, as MCMC chains (or particles) can fail to explore all plausible regions. For emulator methods, the difficulty lies in building a model that can approximate the simulator in all regions of space. For example, stationary covariance functions that assume a constant length-scale throughout space are commonly used in GPs, and may be inappropriate. Other problems arise if $f(\theta)$ varies over too wide a range, which is common if $f(\theta)$ is a likelihood function (Wilkinson, 2014). If we need an emulator of $f(\theta)$ that is valid in all of Θ , then we can look to use a non-stationary covariance function or a more flexible model such as a treed-GP (Gramacy and Lee, 2008). However, for calibration, we only need approximate the simulator when $f(\theta)$ is close to being plausible. In other parts of parameter space, it is only necessary to say θ is implausible with a high degree of confidence. It does not matter if an estimate of $f(\theta)$ is poor, as long as we are correct in saying $f(\theta) \notin \mathcal{P}_C$.

GP predictions are more accurate in regions rich in data. Thus, the key issue when building a GP emulator is the choice of the design points, $\mathcal{D}_\theta = \{\theta_i\}_{i=1}^n$, at which we evaluate the simulator. Space filling designs, such as maximin Latin hypercubes (McKay et al., 2000) or low discrepancy sequences (such as Sobol sequences, Morokoff and Caffisch, 1994) are the default choice of design, and usually lead to reasonable global approximations. But they are less well suited to calibration problems, in which we usually want to focus on just a small region of parameter space.

Instead of a space filling design, we can seek to build the design sequentially: given the current design, we build an emulator that describes our current knowledge of $f(\theta)$. We then use the emulator to decide where next to run the simulator, and so on. The basic idea is as follows:

1. Start with an *a priori* plausible set $\mathcal{P}_\theta^{(0)} = \Theta$.
2. Choose design $\mathcal{D}_\theta^{(1)} = \{\theta_i \in \Theta : i = 1, \dots, n_1\}$, and run the simulator, to get ensemble $\mathcal{D}^{(1)} = \{(\theta_i, C_i = f(\theta_i)) : \theta_i \in \mathcal{D}_\theta^{(1)}\}$.
3. Build emulator $\tilde{f}_{(1)}$ and use it to predict the plausible set $\tilde{\mathcal{P}}_\theta^{(1)}$.
4. Choose new design points $\mathcal{D}_\theta^{(2)}$, and run the simulator to get $\mathcal{D}^{(2)}$.
5. Build emulator $\tilde{f}^{(2)}$ and use it to predict the plausible set $\tilde{\mathcal{P}}_\theta^{(2)}$.
6. Etc.

The details of each step vary in each problem. The plausibility criteria are usually defined so that they become more stringent at each iteration. The first plausibility condition $\mathcal{P}_c^{(1)}$ may be relatively weak, with $\mathcal{P}_c^{(1)}, \mathcal{P}_c^{(2)}, \dots, \mathcal{P}_c^{(W)}$ slowly approaching the final desired criterion $\mathcal{P}_c^{(W)}$. If the difference between $\mathcal{P}_\theta^{(i)}$ and $\mathcal{P}_\theta^{(i+1)}$ is too large, we may find the emulator accuracy is insufficient, causing us to incorrectly rule-out some regions of space (type-I errors). The plausibility criteria can be relaxed by changing the number of measurements we need to match, and the closeness of the required match. Note the superficial similarity to SMC-ABC approaches, in that the approximation is iteratively improved as we learn.

The emulator used at each stage may be based upon all the previous simulator runs, adding new data points in important regions (see below for details), or it can be built from scratch. For example, in Vernon et al. (2010) they build an emulator, $\tilde{f}^{(i)}$, to predict $f(\theta)$ for $\theta \in \mathcal{P}_\theta^{(i-1)}$, the estimated plausible region from the previous iteration. The emulator is not required to predict for $\theta \notin \mathcal{P}_\theta^{(i-1)}$. The benefit of this is that the simulator response is likely to be less variable within $\mathcal{P}_\theta^{(i-1)}$ than in Θ , making it easier to model. The disadvantage is that if some regions are incorrectly ruled to be implausible in iteration $i - 1$, this mistake can never be rectified.

The most important algorithmic decision is the choice of design, $\mathcal{D}_\theta^{(i)}$, at each iteration, i.e., given an emulator, how should we choose locations θ at which to run the simulator? If we only wish the emulator to predict well in $\mathcal{P}_\theta^{(i-1)}$, then we only need a design in $\mathcal{P}_\theta^{(i-1)}$. Vernon et al. (2010) take the approach of seeking to use a space filling design on $\mathcal{P}_\theta^{(i-1)}$, such as a Latin hypercube. To do this, they create a large design on Θ and then reject any point not predicted to lie in $\mathcal{P}_\theta^{(i-1)}$ by $\tilde{f}^{(i)}$, which is also the approach we describe in Section 4. If we instead seek a global emulator valid for all $\theta \in \Theta$, but which is accurate in the important regions, then it can be beneficial to add simulator runs to the design one at a time. The critical regions are those where the emulator is most uncertain about whether $\theta \in \mathcal{P}_\theta$. This is typically either in regions for which we have no data, or near the edge of the plausible region where we are unsure if $\theta \in \mathcal{P}_\theta$ or not given the accuracy of the emulator. If we use a GP emulator, then our prediction of $f(\theta)$ is Gaussian:

$$f(\theta) \sim N(\mu^{(i)}(\theta), \Sigma^{(i)}(\theta))$$

where $\mu^{(i)}$ and $\Sigma^{(i)}$ are the mean and covariance function of $\tilde{f}^{(i)}$. This allows us to calculate the probability that $\theta \in \mathcal{P}_\theta$. For example, if our criterion is

that θ is plausible if $D_- \leq f(\theta) \leq D_+$, then

$$p(\theta) = \mathbb{P}_{\tilde{f}^{(i)}}(\theta \in \mathcal{P}_\theta) = \Phi\left(\frac{D_+ - \mu(\theta)}{\Sigma(\theta)^{\frac{1}{2}}}\right) - \Phi\left(\frac{D_- - \mu(\theta)}{\Sigma(\theta)^{\frac{1}{2}}}\right).$$

In some regions $p(\theta)$ will be close to zero, indicating that we are confident that θ is implausible, and in others close to one, indicating the converse. It is regions in which we are most uncertain, that we wish to target, as these represent parameter values that we can neither rule in nor out. One approach to selecting new design points is to choose points to minimize the entropy of this surface (Hennig and Schuler, 2012; Chevalier et al., 2014). The entropy represents how close to certain knowledge we are. If we let \bar{H} be the average entropy of the emulator prediction of the plausibility surface:

$$\bar{H} = \int -p(\theta) \log p(\theta) - (1 - p(\theta)) \log(1 - p(\theta)) d\theta$$

then we can ask, if we were to add a simulator evaluation at θ , what is the expected value of \bar{H} given the expected resulting information? We can then add θ to the design in order to minimize $\mathbb{E}(\bar{H} | \mathcal{D}^{(i-1)} \cup \{\theta\})$. This approach places new points in regions that most quickly reduce the uncertainty about the plausible region \mathcal{P}_θ .

3.2 A simple climate example

As an illustration of the potential benefit of these techniques, we consider a relatively simple two-box climate simulator (Emanuel, 2002), which models atmospheric and ocean heat transport and storage, with water vapour as a positive feedback. The simulator is useful for the purposes of demonstration, as 10 years of model time takes approximately 5 seconds of CPU time, allowing a large number of model runs to be done. Matlab code for this simulator is available online.

We present the results of a simple history-matching task, calibrating two parameters: DTcrit_conv, the critical vertical temperature gradient that triggers convection, which we allowed to vary in the range [30, 50]; and GAMMA, the emissivity parameter for water vapour, which we varied in the range [1, 2]. We try to find the parameter values that give a global surface temperature between 294.5K and 295.5K once the model is in equilibrium. The CO₂ concentration was set to 560ppm, and all other parameters were set to their default values (EPM, 2010). These choices are arbitrary and only intended for illustration of the methodology.

Applying a simple ABC rejection algorithm and allowing for 1000 simulator evaluations gave us 106 accepted parameter values, which are shown in the left-hand plot in Figure 1, with the red points showing the 10 values accepted after only 100 simulator evaluations. In contrast, the middle and right-hand plots shows the result of using a GP emulator with a maximin Latin hyper-cube (MLH) design of 10 and 30 simulator evaluations. After 10 simulator evaluations, the emulator has some idea of where the plausible region is, but with errors for example in the bottom right hand corner. After 30 simulator evaluations it has accurately found all of the plausible region, with just a little uncertainty at the edge of the region (shown by the grey shading).

This approach, however, relies upon finding a good design. If an accurate emulator results, then it will do well at predicting the plausible region. Here we can see, in the case where we had only 10 design points, that no information is available about the bottom right hand corner of the parameter space, and consequently the model does less well there. As the design is chosen in advance of the simulations being run, finding a good design involves an element of luck.

If we instead use a sequential design and add design points one at a time in order to minimize the expected average entropy of the resulting history-match, then we can significantly improve the speed with which we find \mathcal{P}_θ . The two plots in Figure 2 show the resulting history match after 4 and 10 simulator evaluations. After only 10 simulator evaluations, we have found \mathcal{P}_θ with superior accuracy to that found after 30 simulator evaluations using the MLH design.

Note that the acceptance rate in the ABC algorithm was approximately 10%, considerably higher than in most problems (we had a 1% acceptance rate in the case study described in Section 4). As the acceptance rate decreases, the value of using an emulator to predict \mathcal{P}_θ increases, as the emulator is able to predict where the plausible region is, whereas ABC can only find the region by chance, as it uses no information about the shape of the underlying surface. In contrast, the major advantage of the Monte Carlo approach is that it is less prone to errors (although mixing errors commonly occur in practice), unlike the emulator approach, which can mislead if the fitted model is inaccurate, and thus requires careful supervision.

4 Climate model case study

4.1 The global carbon cycle

Human emissions of CO₂ into the atmosphere are a principal cause of climate change. However, this “anthropogenic” CO₂ does not remain in the atmosphere indefinitely. It is taken up by vegetation and by the oceans, and eventually (after many thousands of years) it is deposited as carbonate sediments at the ocean floor. Understanding these processes is crucial for future climate projections. Climate change is driven by changes in CO₂ concentration and it is therefore determined by the interplay between anthropogenic emissions and the carbon cycle. Many carbon cycle processes are highly uncertain. Projections of year 2100 CO₂ concentrations from different Earth System Models (ESMs) driven by the same assumption of future emissions typically vary by ± 100 ppm (Friedlingstein et al., 2006). This uncertainty range is greater than the total increase to date (2015) due to all historical anthropogenic emissions (~ 120 ppm)

To investigate uncertainties in the global carbon cycle we need a model of appropriate complexity that is capable of resolving the important processes but which is sufficiently numerically efficient. The GENIE-1 intermediate complexity ESM (Holden et al., 2013a) is one such model. The computational speed of GENIE-1 comes mainly from the use of a very simple 2D model of the atmosphere and relatively coarse model resolution (grid cells of $\sim 1,000 \times 1,000$ km). The carbon cycle of GENIE-1 comprises a terrestrial carbon model, a 3D dynamic ocean, dynamic sea ice, ocean biogeochemistry and ocean sediments. Given appropriate model parameter choices, GENIE-1 simulates realistic spatial distributions of carbon storage in vegetation, soil, ocean and carbonate sediment. However, the future response of the climate cycle to ongoing emissions depends upon the specific parameter choices, and will vary even amongst parameter sets that have been constrained to produce similar (and reasonable) modern climate states. To quantify this uncertain response we require an ensemble of simulations that samples widely from plausible input parameter space.

The timescales for different carbon cycle processes vary considerably. Equilibrium timescales are ~ 10 's years for vegetation, ~ 100 's years for soil, $\sim 1,000$'s years for the ocean and $\sim 10,000$'s years for carbonate sediments. In order to simulate an Earth with a carbon cycle in approximate equilibrium (i.e. prior to human interference), a simulation of at least 10,000 years is required¹. Although several orders of magnitude faster than

¹Shorter spin-ups are sufficient for models that neglect sediments.

state-of-the-art ESMs, GENIE-1 requires ~ 4 CPU days to simulate 10,000 real years. The exploration of high-dimensional input space and identification of plausible subspaces is therefore a highly demanding computational problem, which we address through emulator-informed ABC.

4.2 Emulator-informed ABC design

The philosophy of the design approach is to vary key model parameters over the entire range of plausible values and to accept those parameter combinations that lead to climate states that cannot be uncontroversially ruled out as implausible (Edwards et al., 2011). We are seeking to explore all plausible simulator realisations in order to capture the range of possible feedback strengths. The input ranges we apply, Θ , are generally broader than ranges that are applied in model tuning exercises. This is in part to enable us to fully quantify model behaviour over plausible parameter space, \mathcal{P}_θ , and in part to improve the validity of the ensemble for application to diverse climate states, such as the Last Glacial Maximum.

The experimental set-up is described in Holden et al. (2013b). We varied 24 model parameters in the ensemble. The choice of parameters was governed by consideration of the processes that are thought to contribute to the natural variability of atmospheric CO₂ on glacial-interglacial timescales (Kohfeld and Ridgwell, 2009) and hence to which the distribution of carbon may be sensitive in general. Five atmospheric parameters were varied. These parameters control the spatial distribution of simulated temperature and precipitation, and hence drive changes in vegetation, sea-ice coverage and ocean circulation. Five parameters were varied in the vegetation model, controlling photosynthesis and respiration rates. Five ocean parameters were varied. These control ocean circulation, and hence the spatial distribution of carbon, alkalinity, dissolved oxygen and nutrients in the ocean. Sea-ice diffusivity was varied, primarily because of its effect on ocean circulation by altering the transport of freshwater. Nine ocean biogeochemistry parameters were varied. These parameters drive changes in the rates of atmosphere-ocean gas exchange, plankton photosynthesis and the remineralisation of the organic products of this photosynthesis. The rate of remineralisation controls the transport of carbon from the surface of the ocean to the deep.

A 500-member ensemble of 25,000-year simulations was first performed using a maximin Latin hypercube (MLH) design². The plausibility of each

²We note that while efficiencies can be gained in certain applications by initialising each ensemble member with output from an existing equilibrium simulation, such an approach is not likely to be useful here as our approach is designed to sample widely differing Earth system states.

simulator run was evaluated using eight different output quantities, usually termed *metrics* in the climate literature. These simple metrics impose no constraints on the spatial distribution of modelled outputs. They instead provide global-scale constraints on atmosphere (global average temperature), ocean (strength of North Atlantic overturning and Antarctic Deep Water formation), Antarctic sea-ice coverage, global vegetation carbon, global soil carbon, ocean biogeochemistry (average dissolved oxygen concentration in the global ocean) and ocean sediments (the average percentage of CaCO_3 in the surface sediment). Only four of the 500 MLH simulations were found to satisfy all eight plausibility constraints, which given that $\dim \theta = 24$, is insufficient for any meaningful statistical analysis. The MLH ensemble took more than ten years of computing to complete, demonstrating that a naive application of ABC is infeasible for this application.

As described in Section 3, we can use emulators to guide the search to find plausible regions of parameter space. Regression-based emulators, including linear and quadratic terms, were built for each of the eight metrics (outputs) specified above. Prior to fitting, variables were linearly mapped onto the range $[-1,1]$ so that odd and even terms were orthogonal, aiding variable selection. The models were built using a stepwise model selection scheme, initially using the Akaike Information Criterion as the selection criterion, and then subsequently shrunk further by applying the more stringent Bayes Information Criterion. This procedure of first growing the model beyond the BIC constraint and then shrinking helps to avoid local minima in the stepwise search.

Parameters were then sampled uniformly from the a prior plausible region and the emulators used to predict if they would lead to plausible simulations. Parameters were accepted as potentially plausible when the emulators predicted plausible values for all eight metrics. The plausibility ranges used were based on the observed climate record, the simulator discrepancy, and the emulator accuracy. Each accepted parameter combination was then used as a design point in a further simulation.

As simulations completed, the emulators were rebuilt using the additionally available data. This process progressively improved the success rate of the emulator predictions (i.e. the percentage of emulator predicted plausible parameters that led to plausible simulations) from 24% to 65%. In total, the simulator was run for 1,000 parameter values predicted to be plausible by the emulator. This produced 885 completed simulations of which 471 were plausible (the remaining 115 simulations terminated before completion, a common occurrence with climate simulators). This 471-member plausible set forms the Emulator Filtered Plausibility Constrained (EFPC) ensemble. The generation of these simulations required a further 25 years of computing

time. Without the ~ 50 -fold increase in efficiency gained by using an emulator to predict the plausible region, this would have required an infeasible amount (more than 1000 years) of CPU time.

While it is clear that ABC strongly constrains the outputs (metrics) that are explicitly filtered for, it is worth noting that it indirectly constrains all aspects of the Earth system and leads to improved simulated magnitudes *and spatial distributions* of state variables generally. Figure 3 provides an illustrative output of the EFPC ensemble, and of the benefits of the ABC filtering. The figure illustrates cross-sections of ocean alkalinity through the Atlantic and Pacific oceans, comparing ensemble means of the unfiltered MLH simulations (left) and filtered EFPC simulations (centre) with observational data (right). Ocean alkalinity exerts a strong control on atmospheric carbon dioxide by determining the degree to which dissolved carbon dioxide is dissociated into bicarbonate and carbonate ions, in turn determining the rates of exchange of dissolved carbon in the ocean with the atmosphere (carbon dioxide) and the sediments (calcium carbonate). Alkalinity is not directly constrained by the ABC metrics, but its distribution is influenced by them, for instance through the constraints imposed on ocean circulation strength and the sediment carbonate concentration. Relative to the MLH ensemble, the EFPC ensemble shows elevated surface concentrations, decreased concentrations in the deep Atlantic (apparently associated with the Atlantic overturning circulation in the unfiltered ensemble) and increased penetration of high alkalinity towards Southern latitudes in the deep Pacific. Although discrepancies with observations remain, which may reflect structural deficiencies in the simulator, each of these trends produces better ensemble-averaged agreement with the observed distribution.

4.3 Applications

Although the use of ABC to derive a posterior distribution is useful in itself, our primary motivation is to identify a set of plausible parameters for application to diverse simulation problems. The EFPC parameter set has been used in a range of experiments, considering both past and future climate change. For clarity, it is worth emphasising that while these experiments did not use ABC directly, they were all rendered tractable by the use of emulator-informed ABC to design the underlying simulation ensemble. A selection of these experiments are summarised below, each with a focus on a different category of application.

4.3.1 Probabilistic simulation outputs: the uncertain response of the carbon cycle to anthropogenic CO₂ emissions

We (PBH and NRE) contributed a suite of carbon cycle experiments for the Fifth Assessment Report (AR5) of the Intergovernmental Report on Climate Change (IPCC). Fifteen Intermediate Complexity ESMs from around the world performed these experiments. The focus was on historical change (Eby et al., 2013) and long-term future change (Zickfeld et al., 2013), considering long timescale problems that are not tractable by state-of-the-art ESMs, thus requiring the use of reduced complexity models such as GENIE-1. Forty seven experiments were performed.

We applied a subset of the EFPC parameter set, in part to aid computational tractability, in view of the 47 separate experiments required, and in part to eliminate a bias in the transient response of the ensemble. The EFPC parameter set is constrained to simulate a plausible preindustrial climate, but no constraint was imposed upon the dynamic response to anthropogenic emissions. Four important model parameters were not constrained by preindustrial plausibility, two relating to cultivated vegetation (deforestation for agriculture was neglected in the preindustrial simulations), a parameter controlling the direct effect of CO₂ on photosynthesis (“CO₂ fertilisation”, see following section) and a parameter controlling the uncertain effect of clouds on the Earth’s radiation budget in a warmer planet. The dynamic response was therefore filtered through a historical forcing experiment, which imposed anthropogenic forcing, including CO₂ emissions, since preindustrial times in an EFPC ensemble of transient simulations. Twenty parameter sets, selected at random from the EFPC parameter sets, but constrained to approximately reproduce the present day atmospheric CO₂ concentration, were accepted and applied to the IPCC experiments.

We do not attempt to summarise the results of these extensive multi-model comparisons here, but note that the GENIE-1 perturbed-parameter ensemble was found to provide an unbiased representation of the multi-model ensemble, being approximately centred on the mean of the fifteen models and with comparable uncertainty. These uncertainties were presented in a related model intercomparison paper (Joos et al., 2013).

4.3.2 Calibrating model parameters: the strength of the terrestrial carbon sink

The IPCC experiments revealed a general tendency of intermediate complexity ESMs to understate the magnitude of the terrestrial carbon sink (the anthropogenic CO₂ taken up by vegetation on land). The major uncertainty

in the terrestrial sink relates to CO₂ fertilisation. Experimental evidence almost without exception shows a stimulation of leaf photosynthesis when plants are exposed to elevated CO₂ (Körner, 2006). In addition to this direct affect on photosynthesis, the short time-scale physiological effect of reduced stomatal opening increases water-use efficiency and additionally increases the efficiency of photosynthesis (Field et al., 1995). However, the strength of the fertilisation effect is poorly quantified, especially under natural conditions. Some studies have failed to detect a measurable effect in nature, while others suggest that any effects may be short term, as CO₂ is only one of a number of potentially limiting factors on plant growth (Körner, 2006).

We addressed this calibration problem in Holden et al. (2013a). Using output from a 671-member ensemble of transient GENIE-1 simulations derived from the EFPC parameter sets we built an emulator of the change in atmospheric CO₂ concentration change since the preindustrial period. We then applied this emulator to sample the 28-dimensional input parameter space. A Bayesian calibration suggests that the increase in gross primary productivity (GPP) in response to a doubling of CO₂ from preindustrial values is very likely (90% confidence) to exceed 20%, with a most likely value of 40-60%.

4.3.3 Model understanding: what determines the spatial distribution of dissolved carbon in the ocean?

In Holden et al. (2013b) we applied the EFPC ensemble to a transient experiment over the recent industrial era (1858 to 2008 AD). The temporal evolution of atmospheric CO₂ and its isotopic composition are known from observational data, and these simulated quantities were made to follow the observations through a relaxation term. The objective of the experiment was to better understand the mechanisms by which the anthropogenic CO₂ emissions are taken up by the ocean.

To achieve this, we analysed the change in distributions of ocean concentrations of dissolved inorganic carbon (DIC) and its stable isotope $\delta^{13}C$, considering two-dimensional latitudinal-vertical transects through the Atlantic and Pacific. These two transects were combined into a single vector for each simulation (to ensure inter-basin effects were consistently represented), and these vectors were combined into an ensemble matrix. Singular vector decomposition was applied to the DIC and $\delta^{13}C$ matrices in order to extract the dominant modes of their spatial variability across the ensemble. Emulators of the component scores elicited further understanding of these modes by identifying which model parameters were driving each mode of variability.

This, together with physical interpretation of the spatial patterns of each

mode, enabled us to identify the principal processes driving them, on the assumption that the dominant parameters governing uncertainty in the response of each mode could be identified with the most important parameterised processes controlling the respective modes. We showed that the main processes governing the uptake of anthropogenic CO_2 and $\delta^{13}\text{C}$ are quite distinct: an important conclusion because observations of the isotopic composition are used to infer rates of ocean CO_2 uptake. Uncertainty in anthropogenic $\delta^{13}\text{C}$ uptake is dominated by air-sea gas exchange, which explains 63% of modelled variance. This mode of variability is largely absent from the ensemble variability in CO_2 uptake, which is instead driven by uncertainties in mixing rates between the surface and deep oceans.

4.3.4 Coupling applications: coupling climate models and climate change impact models

The evaluation of climate impacts requires coupling climate models, impact models and economic models together within an “Integrated Assessment Model” (IAM) framework. In such couplings, climate data (e.g. regional temperature, precipitation) are passed to the IAM for computation of climate impact functions, and the IAM passes back anthropogenic forcing (such as CO_2 emissions or land use change). Computational demands mean that it is generally infeasible to couple complex climate models into IAMs. Various approaches are taken to address this, using either simplified models or statistical representations of more complex models. Recently, effort has focussed on the use of emulators of climate models as surrogates for the simulator in these coupling applications.

Economic models provide projections of CO_2 emissions. They typically convert emissions into concentrations through the use of simple “box-models”, describing rates of carbon transfer between the atmospheric, terrestrial and oceanic reservoirs. We have recently applied the EFPC parameter set to build an emulator of the GENIE-1 carbon cycle model for incorporation into integrated assessment models (Foley et al., in prep). An 86-member subset of the EFPC parameter set was used to generate an ensemble of future climate-carbon cycle experiments, with future emissions prescribed as modified Chebyshev polynomials.

The emulation approach followed the “1-step” dimensionally-reduced emulation methodology of Holden et al. (in press), emulating a singular value decomposition of the ensemble outputs. Emulators of the first four component scores were derived as functions of the 28 model parameters and the 6 concentration profile coefficients. The emulator outputs are, unsurprisingly, dominated by the Chebyshev coefficients. However, uncertainty

for a given forcing scenario is generated through emulator dependencies on GENIE-1 parameters. The resulting carbon cycle emulator has been coupled into an integrated assessment framework that also includes a macroeconomic model of the global economy E3MG (Mercure et al., 2014), an agent-based model of technology substitution dynamics FTT-power (Mercure, 2012) and a spatiotemporally resolved emulator of the climate system (Holden et al., 2014). We have applied the framework to assess the impact on the climate of emissions reduction policies in the electricity sector (Mercure et al., 2014), addressing the cascade of uncertainty through the coupled system.

5 Future applications

It may never be possible to apply statistical approaches to robustly calibrate a truly state-of-the-art climate simulator. They are defined by the limits of available computing power, and consequently very few simulations are possible with these models. This begs the important question of how far could one go with simulator complexity and still be able apply these methods. We have demonstrated the application of emulator-informed ABC to generate a 471-member ensemble of a model that takes ~ 10 days to perform each simulation. The computational constraints ultimately determined the number of parameters we could vary; a rule of thumb dictates that we use a minimum of 10 ensemble members for each varied active input (Loeppky et al., 2009). It is worth noting that a useful ensemble varying only, say, 5 parameters would need ~ 50 simulations, and could have been achieved for a 10-fold slower model.

The improvements in methodology demonstrated in Section 3.2, suggest efficiencies that should significantly extend the applicability of the approach. The use of GP emulation generally allows a better statistical model than linear regression, and therefore would be expected to improve the success rate of the emulator filtering. This will certainly be the case when a parametric mean function is used and the GP is applied only to emulate the residual. The uncertainty estimates provided by the GP should also improve the success rate of the emulator filtering, for instance, by only accepting parameters for which there is a high probability of plausibility. Furthermore, a significant improvement arises from the use of a sequential design process, which was shown to yield a 3-fold increase in efficiency in our example. For more complex simulators, we will want to make use of parallel computation. The sequential approach then changes from adding one design point at a time, to adding d , where d is the number of available cores. Finding the d optimal points that minimize the expected entropy is difficult, and is an area of ac-

tive research, but even suboptimal designs can give significant improvements over the default space-fillings designs. For stochastic simulators, many of the same techniques can be applied. The likelihood function now needs to be estimated, significantly increasing the difficulty, but progress is being made in this direction (Oakley and Youngman, 2014; Meeds and Welling, 2014; Wilkinson, 2014; Gutmann and Corander, 2015).

These improvements in efficiency should render application to “previous-generation” ESMs such as HadCM3³ tractable on multi-node computing clusters, certainly so on distributed computing systems such as climateprediction.net, which last year facilitated more than 7,500 years of climate modelling on the personal computers of the general public.

References

- A box climate model: EPcm. model documentation v4. www.sp.ph.ic.ac.uk/~aczaja/EP_ClimateModel.html, 2010.
- I. Andrianakis, I. R. Vernon, N. McCreesh, T. J. McKinley, J. E. Oakley, R. N. Nsubuga, M. Goldstein, and R. G. White. Bayesian history matching of complex infectious disease models using emulation: A tutorial and a case study on hiv in uganda. *PLoS computational biology*, 11(1):e1003968, 2015.
- J. Annan, J. Hargreaves, N. Edwards, and R. Marsh. Parameter estimation in an intermediate complexity earth system model using an ensemble Kalman filter. *Ocean modelling*, 8(1):135–154, 2005.
- G. Biau, F. Cérou, A. Guyader, et al. New insights into Approximate Bayesian Computation. In *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*, volume 51, pages 376–403. Institut Henri Poincaré, 2015.
- M. G. Blum, M. A. Nunes, D. Prangle, S. A. Sisson, et al. A comparative review of dimension reduction methods in approximate Bayesian computation. *Statistical Science*, 28(2):189–208, 2013.
- J. Brynjarsdóttir and A. O’Hagan. Learning about physical parameters: The importance of model discrepancy. *Inverse Problems*, 30(11):114007, 2014.
- S. Castruccio, D. J. McInerney, M. L. Stein, F. Liu Crouch, R. L. Jacob, and E. J. Moyer. Statistical emulation of climate model projections based on precomputed GCM runs*. *Journal of Climate*, 27(5):1829–1844, 2014.

³HadCM3 performs more than 10 years per CPU day on eight nodes of a linux cluster.

- C. Chevalier, D. Ginsbourger, J. Bect, E. Vazquez, V. Picheny, and Y. Richet. Fast parallel kriging-based stepwise uncertainty reduction with application to the identification of an excursion set. *Technometrics*, 56(4):455–465, 2014.
- P. S. Craig, M. Goldstein, A. H. Seheult, and J. A. Smith. Pressure matching for hydrocarbon reservoirs: a case study in the use of Bayes linear strategies for large computer experiments. In *Case studies in Bayesian statistics*, pages 37–93. Springer, 1997.
- M. Eby, A. J. Weaver, K. Alexander, K. Zickfeld, A. Abe-Ouchi, A. Cimadoribus, E. Cresspin, S. Drijfhout, N. Edwards, A. Eliseev, et al. Historical and idealized climate model experiments: an intercomparison of Earth system models of intermediate complexity. *Climate of the Past*, 9:1111–1140, 2013.
- N. R. Edwards, D. Cameron, and J. Rougier. Precalibrating an intermediate complexity climate model. *Climate dynamics*, 37(7-8):1469–1482, 2011.
- K. Emanuel. A simple model of multiple climate regimes. *Journal of Geophysical Research*, 107(D9):ACL–4, 2002.
- C. Field, R. Jackson, and H. Mooney. Stomatal responses to increased CO₂: implications from the plant to the global scale. *Plant, Cell & Environment*, 18(10):1214–1225, 1995.
- G. Flato, J. Marotzke, et al. Evaluation of climate models. *Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*, pages 741–866, 2013.
- A. Foley, J.-F. Mercure, P. Salas, P. Holden, N. Edwards, H. Pollitt, and U. Chewpreecha. Modelling the climate impacts of mitigation policies in the energy sector: An integrated assessment framework using climate model emulation. *for Earth System Dynamics*, in prep.
- P. Friedlingstein, P. Cox, R. Betts, L. Bopp, W. Von Bloh, V. Brovkin, P. Cadule, S. Doney, M. Eby, I. Fung, et al. Climate-carbon cycle feedback analysis: Results from the C4MIP model intercomparison. *Journal of Climate*, 19(14):3337–3353, 2006.
- R. B. Gramacy and H. K. Lee. Bayesian treed gaussian process models with an application to computer modeling. *Journal of the American Statistical Association*, 103(483), 2008.
- M. U. Gutmann and J. Corander. Bayesian optimization for likelihood-free inference of simulator-based statistical models. *arXiv preprint arXiv:1501.03291*, 2015.

- G. R. Harris, D. M. Sexton, B. B. Booth, M. Collins, and J. M. Murphy. Probabilistic projections of transient climate change. *Climate Dynamics*, 40:2937–2972, 2013.
- P. Hennig and C. J. Schuler. Entropy search for information-efficient global optimization. *The Journal of Machine Learning Research*, 13(1):1809–1837, 2012.
- D. Higdon, J. Gattiker, B. Williams, and M. Rightley. Computer model calibration using high-dimensional output. *Journal of the American Statistical Association*, 103(482), 2008.
- P. Holden and N. Edwards. Dimensionally reduced emulation of an AOGCM for application to integrated assessment modelling. *Geophysical Research Letters*, 37(21), 2010.
- P. Holden, N. Edwards, D. Gerten, and S. Schaphoff. A model-based constraint on CO₂ fertilisation. *Biogeosciences*, 10(1):339–355, 2013a.
- P. Holden, N. Edwards, S. Müller, K. Oliver, R. Death, and A. Ridgwell. Controls on the spatial distribution of oceanic $\delta^{13}C_{DIC}$. *Biogeosciences*, 10:1815–1833, 2013b.
- P. Holden, N. Edwards, P. Garthwaite, K. Fraedrich, F. Lunkeit, E. Kirk, M. Labriet, A. Kanudia, and F. Babonneau. PLASIM-ENTSem v1. 0: a spatio-temporal emulator of future climate change for impacts assessment. *Geoscientific Model Development*, 7(1):433–451, 2014.
- P. Holden, N. Edwards, P. Garthwaite, and R. Wilkinson. Emulation and interpretation of high-dimensional climate model output. *Journal of Applied Statistics*, in press.
- P. B. Holden, N. Edwards, K. Oliver, T. Lenton, and R. Wilkinson. A probabilistic calibration of climate sensitivity and terrestrial carbon change in GENIE-1. *Climate Dynamics*, 35(5):785–806, 2010.
- J. L. Holloway Jr and S. Manabe. Simulation of climate by a global general circulation model: I. hydrologic cycle and heat balance. *Monthly Weather Review*, 99(5):335–370, 1971.
- F. Joos, R. Roth, J. Fuglestad, G. Peters, I. Enting, W. v. Bloh, V. Brovkin, E. Burke, M. Eby, N. Edwards, et al. Carbon dioxide and climate impulse response functions for the computation of greenhouse gas metrics: a multi-model analysis. *Atmospheric Chemistry and Physics*, 13(5):2793–2825, 2013.
- M. Kennedy and A. O’Hagan. Bayesian calibration of computer models (with discussion). *J. R. Statist. Soc. Ser. B*, 63:425–464, 2001.

- K. E. Kohfeld and A. Ridgwell. Glacial-Interglacial Variability in Atmospheric CO₂. *Surface ocean-lower atmosphere processes*, pages 251–286, 2009.
- C. Körner. Plant CO₂ responses: an issue of definition, time and resource supply. *New phytologist*, 172(3):393–411, 2006.
- L. Lee, K. Carslaw, K. Pringle, and G. Mann. Mapping the uncertainty in global CCN using emulation. *Atmospheric Chemistry and Physics*, 12(20):9739–9751, 2012.
- J. L. Loeppky, J. Sacks, and W. J. Welch. Choosing the sample size of a computer experiment: A practical guide. *Technometrics*, 51(4), 2009.
- J. Marquis, Y. Richardson, P. Markowski, D. Dowell, J. Wurman, K. Kosiba, P. Robinson, and G. Romine. An investigation of the Goshen County, Wyoming, tornadic supercell of 5 June 2009 using EnKF assimilation of mobile mesonet and radar observations collected during VORTEX2. Part I: Experiment design and verification of the EnKF analyses. *Monthly Weather Review*, 142(2):530–554, 2014.
- T. Mauritsen, B. Stevens, E. Roeckner, T. Crueger, M. Esch, M. Giorgetta, H. Haak, J. Jungclaus, D. Klocke, D. Matei, et al. Tuning the climate of a global model. *Journal of Advances in Modeling Earth Systems*, 4(3), 2012.
- M. D. McKay, R. J. Beckman, and W. J. Conover. A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics*, 42(1):55–61, 2000.
- E. Meeds and M. Welling. GPS-ABC: Gaussian process surrogate approximate Bayesian computation. *arXiv preprint arXiv:1401.2838*, 2014.
- J.-F. Mercure. FTT: Power: A global model of the power sector with induced technological change and natural resource depletion. *Energy Policy*, 48:799–811, 2012.
- J.-F. Mercure, H. Pollitt, U. Chewpreecha, P. Salas, A. M. Foley, P. B. Holden, and N. R. Edwards. The dynamics of technology diffusion and the impacts of climate policy instruments in the decarbonisation of the global electricity sector. *Energy Policy*, 73:686–700, 2014.
- W. J. Morokoff and R. E. Caflisch. Quasi-random sequences and their discrepancies. *SIAM J. Sci. Comput.*, 15:1251–1279, 1994.
- J. M. Murphy, D. M. Sexton, D. N. Barnett, G. S. Jones, M. J. Webb, M. Collins, and D. A. Stainforth. Quantification of modelling uncertainties in a large ensemble of climate change simulations. *Nature*, 430(7001):768–772, 2004.

- J. E. Oakley and B. D. Youngman. Calibration of complex computer simulators using likelihood emulation. *arXiv preprint arXiv:1403.5196*, 2014.
- A. O’Hagan. Bayesian analysis of computer code outputs: a tutorial. *Reliability Engineering & System Safety*, 91(10):1290–1300, 2006.
- R. Olson, R. Sriver, M. Goes, N. M. Urban, H. D. Matthews, M. Haran, and K. Keller. A climate sensitivity estimate using Bayesian fusion of instrumental observations and an Earth System model. *Journal of Geophysical Research: Atmospheres (1984–2012)*, 117(D4), 2012.
- O. K. Oyebamiji, N. R. Edwards, P. B. Holden, P. H. Garthwaite, S. Schaphoff, and D. Gerten. Emulating global climate change impacts on crop yields. *Statistical Modelling*, page 1471082X14568248, 2015.
- J. Rougier. Probabilistic inference for future climate using an ensemble of climate model evaluations. *Climatic Change*, 81(3-4):247–264, 2007.
- J. Rougier and M. Goldstein. Climate simulators and climate projections. *Annual Review of Statistics and Its Application*, 1:103–123, 2014.
- J. Rougier, D. M. Sexton, J. M. Murphy, and D. Stainforth. Analyzing the climate sensitivity of the HadSM3 climate model using ensembles from different but related experiments. *Journal of Climate*, 22(13):3540–3557, 2009.
- J. Sacks, W. J. Welch, T. J. Mitchell, and H. P. Wynn. Design and analysis of computer experiments. *Stat. sci.*, 4:409–423, 1989.
- B. Sansó, C. E. Forest, D. Zantedeschi, et al. Inferring climate system properties using a computer model. *Bayesian Analysis*, 3(1):1–37, 2008.
- T. J. Santner, B. J. Williams, and W. I. Notz. *The design and analysis of computer experiments*. Springer Verlag, 2003.
- K. Sham Bhat, M. Haran, R. Olson, and K. Keller. Inferring likelihoods and climate system characteristics from climate models and multiple tracers. *Environmetrics*, 23(4):345–362, 2012.
- I. Vernon, M. Goldstein, and R. G. Bower. Galaxy formation: a Bayesian uncertainty analysis. *Bayesian Analysis*, 5:619–669, 2010.
- R. D. Wilkinson. Bayesian calibration of expensive multivariate computer experiments. *Large-Scale Inverse Problems and Quantification of Uncertainty*, pages 195–215, 2010.
- R. D. Wilkinson. Approximate Bayesian computation (ABC) gives exact results under the assumption of model error. *Stat. Appl. Genet. Mo. B.*, 12:129–141, 2013.

- R. D. Wilkinson. Accelerating ABC methods using Gaussian processes. *JMLR Workshop and Conference Proceedings Volume 33: Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics*, 33:1015–1023, 2014.
- D. Williamson, M. Goldstein, L. Allison, A. Blaker, P. Challenor, L. Jackson, and K. Yamazaki. History matching for exploring and reducing climate model parameter space using observations and a large perturbed physics ensemble. *Climate dynamics*, 41(7-8):1703–1729, 2013.
- K. Zickfeld, M. Eby, A. J. Weaver, K. Alexander, E. Cresspin, N. R. Edwards, A. V. Eliseev, G. Feulner, T. Fichefet, C. E. Forest, et al. Long-term climate change commitment and reversibility: An EMIC intercomparison. *Journal of Climate*, 26(16):5782–5809, 2013.

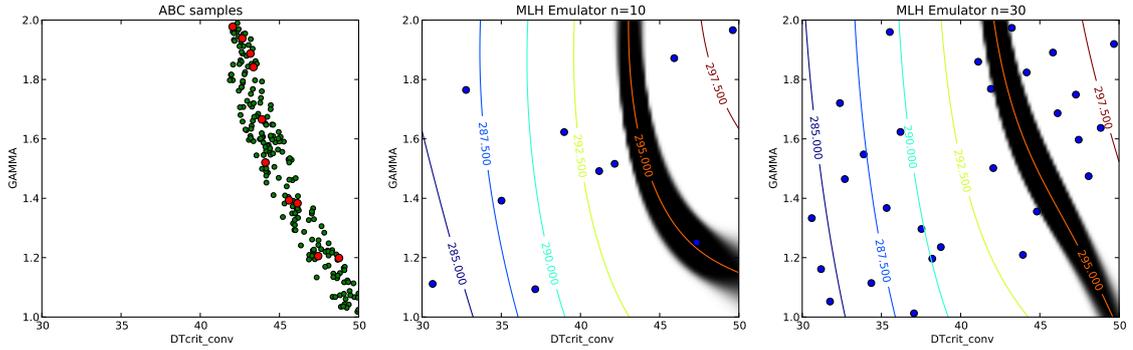


Figure 1: Left: Accepted samples from the rejection ABC algorithm after 100 (red) and 1000 (green) simulator evaluations. Middle and right: The estimated plausible region using an emulator trained with a maximin Latin hypercube design (points shown in blue) with 10 (middle) and 30 (right) simulator evaluations. The shading indicates the estimated value of $\mathbb{P}(\theta \in \mathcal{P}_\theta)$. The contour lines are the estimated response surface $f(\theta)$.

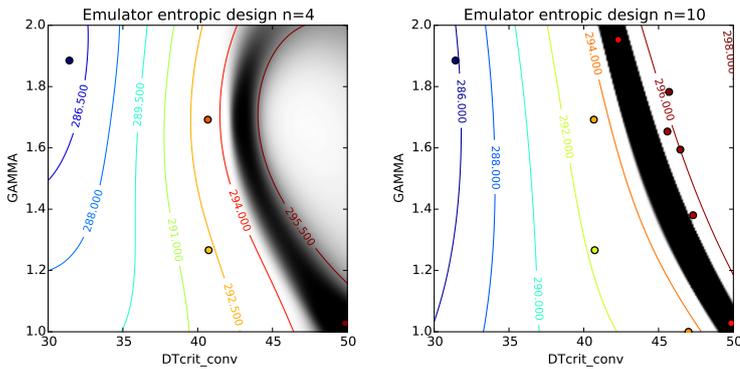


Figure 2: Results from using an entropy based sequential design. The left-hand column shows the estimated response surface (contours) and $\mathbb{P}(\theta \in \mathcal{P}_\theta)$ (shading), with the design points overlaid. The large red point is the most recently added point. The right-hand column shows the entropy surface. The top row uses four simulator evaluations, and the bottom row uses 10 simulator evaluations, all added according to the entropy criterion.

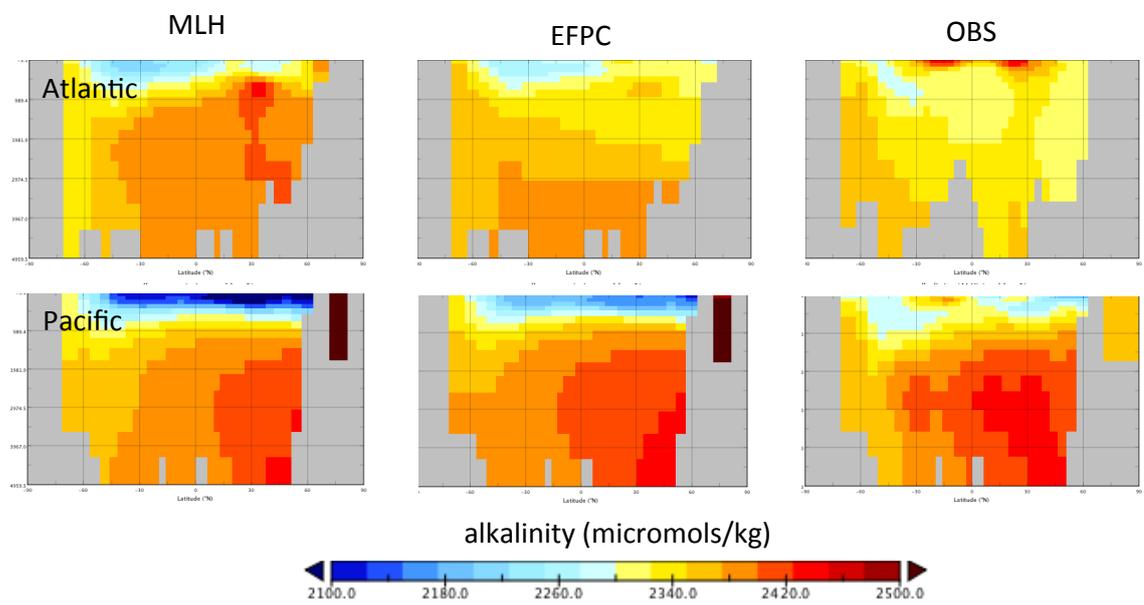


Figure 3: Cross-sections of ocean alkalinity through the Atlantic (25°W) and Pacific (155°W) oceans. The figure compares the mean of the training MLH ensemble (left panels) and the plausibility filtered EFPC ensemble (centre panels) with observations (right panels).