

1 **Identifying and removing structural biases in climate**
2 **models with history matching**

3 **Daniel Williamson · Adam Blaker ·**
4 **Charlotte Hampton · James Salter**

5
6 Received: date / Accepted: date

7 **Abstract** We describe the method of history matching, a method currently used
8 to help quantify parametric uncertainty in climate models, and argue for its use in
9 identifying and removing structural biases in climate models at the model devel-
10 opment stage. We illustrate the method using an investigation of the potential to
11 improve upon known ocean circulation biases in a coupled non-flux-adjusted cli-
12 mate model (the third Hadley Centre Climate Model; HadCM3). In particular, we
13 use history matching to investigate whether or not the behaviour of the Antarc-
14 tic Circumpolar Current (ACC), which is known to be too strong in HadCM3,
15 represents a structural bias that could be corrected using the model parameters.
16 We find that it is possible to improve the ACC strength using the parameters and
17 observe that doing this leads to more realistic representations of the sub-polar
18 and sub-tropical gyres, sea surface salinities (both globally and in the North At-
19 lantic), sea surface temperatures in the sinking regions in the North Atlantic and
20 in the Southern Ocean, North Atlantic Deep Water flows, global precipitation,
21 wind fields and sea level pressure. We then use history matching to locate a region
22 of parameter space predicted not to contain structural biases for ACC and SSTs
23 that is around 1% of the original parameter space. We explore qualitative features
24 of this space and show that certain key ocean and atmosphere parameters must be
25 tuned carefully together in order to locate climates that satisfy our chosen met-
26 rics. Our study shows that attempts to tune climate model parameters that vary
27 only a handful of parameters relevant to a given process at a time will not be as
28 successful or as efficient as history matching.

29 **Keywords** Tuning · Ensembles · Emulators · HadCM3 · Climate Model

D. Williamson · J. Salter
College of Engineering, Mathematics and Physical Sciences, University of Exeter, Exeter, UK
E-mail: d.williamson@exeter.ac.uk

A. T. Blaker
National Oceanography Centre, Southampton, UK, SO14 3ZH

1 Introduction

One of the principal challenges facing the climate modelling community is the removal of systematic or structural errors in generalised circulation models (GCMs) (Randall et al, 2007). So called “known biases” in a GCM drive the development and improvement of these models. For example, a motivation for the development of HadGEM2 was the improvement of the model ENSO compared with its predecessor HadGEM1 (Martin et al., 2010). The Hadley Centre models are not alone in this regard, with each new version of a group’s GCM containing biases that the modellers speculate can be improved or removed with better parameterization schemes or finer resolution (for example, Watanabe et al., 2010; Gent et al., 2011).

As the desire to remove these structural errors drives increases in model resolution and development of new code and parameterization schemes, it is important to know that the errors in question really do represent structural deficiencies of the model and are not merely an artefact of poor tuning of the current parameterization schemes. GCMs necessarily consist of many parameterized schemes designed to approximate the physics in the real world on the grid scale of the model. Each scheme contains a number of parameters whose value must be fixed in order to run the climate model. A major part of the development of a new climate model represents the tuning of these parameters and schemes in order to ensure the resulting model climate is consistent with observations over a number of chosen metrics.

A climate model bias represents a structural error if that bias cannot be removed by changing the parameters without introducing more serious biases to the model. Hence, to state that a climate model bias represents a structural error is to assume that the model has been optimally tuned and yet fails to adequately represent the metric in question. In this paper we argue that any GCM is highly unlikely to be optimally tuned due to the way the parameters are usually selected by modellers.

Climate models are usually tuned on a process by process basis and by trial and error (Severijns and Hazeleger, 2005). An individual process or module is selected and one or two parameters thought to drive that process are changed. If the change moves the process closer to reality, the change is accepted. For example, Acreman and Jeffery (2007) change two parameters in the Kraus and Turner (1967) mixed layer scheme. The study was used to fix these parameters in studies using the UK Met Office’s 3rd Hadley centre model HadCM3 (Gordon et al., 2000; Pope et al., 2000; Collins et al., 2007). Martin et al. (2011) reduce background tracer diffusivity in the ocean by an order of magnitude to improve SST profiles for HadGEM2. There are very few guidelines for tuning the parameters of a climate model. The 4th IPCC report (Solomon et al., 2007, sect. 8.1.3) gives 2 guidelines for tuning. The first is that observation-based constraints on parameter ranges should not be exceeded. The second is that climate model performance is only judged with respect to observational constraints not used in tuning. Whilst the second of these seems sensible, the first is questionable (why is it necessarily true that a numerically integrated solution to climate model equations over a relatively coarse spatial grid be most informative for the true climate when theoretically observable parameters are within their real-world ranges?). Taken together, these guidelines offer little actual instruction for tuning.

Mauritsen et al. (2012) offer a tuning protocol, which was used to develop the latest version of the MPI-ESM model. Their protocol is based on identifying

78 model biases and targeting those in particular by iteration through steps that
79 first involve short runs with prescribed SSTs to find promising parameter choices.
80 These choices are then subjected to longer simulations and compared to observed
81 climate. If they are still promising, they are changed in the coupled model and the
82 resulting model climate is evaluated.

83 Though the use of short runs in preliminary tuning steps seems promising, the
84 focus on tuning only parameters influencing specific processes using an uncoupled
85 version of the model is problematic. Experiments such as these represent “one
86 factor at a time” (OFAT) designs. The hope is that after changing the param-
87 eter choices individually in order to improve each process, the coupled climate
88 model will also have improved. This type of experimental design is well known
89 in the statistics literature for being both inefficient and dangerous (Fisher, 1926;
90 Friedman and Savage, 1947; Daniel, 1973). In particular, if parameters controlling
91 different processes interact, that is, if by changing them simultaneously in some
92 way the effect is different to changing them separately, OFAT type designs cannot
93 find these interactions. Partly because of this and partly due to inefficiency, this
94 type of design is prone to missing optimal settings of the parameters.

95 More formal procedures for climate model tuning in the literature do exist.
96 For example, Severijns and Hazeleger (2005) treat tuning as a global optimization
97 problem which they solve using the downhill simplex method (a numerical min-
98 imisation algorithm). A class of data assimilation methods approach tuning with
99 respect to the key uncertainties: observation error and structural error. These
100 methods, for example, based on the ensemble Kalman filter, combine the param-
101 eters with the climate model state vector in order to fine tune a model, and have
102 been applied to intermediate complexity climate models (Annan et al., 2005c; Har-
103 greaves et al., 2004) and to the atmosphere-only component of a GCM (Annan et
104 al., 2005a,b). However, as yet, they have not been applied to successfully tune a
105 coupled atmosphere-ocean GCM and, Rougier (2013) states that there is reason
106 to think that this type of tuning method is intractable in the full parameter space.

107 Another problem with data assimilation approaches to parameter tuning is
108 that the parameters, and hence the model physics, are allowed to vary in time.
109 This means that the final parameter choice, that which corresponds to the value of
110 the augmented state vector following assimilation of the most recent observations,
111 need not represent a model that would reproduce an acceptable solution to the
112 underlying model equations over the full assimilation period. In fact the model
113 is constantly tuned so that the solution is not dynamically consistent. We might
114 think of these solutions as representing worlds with “transient physics”, which can
115 be seen as undermining the key assumptions made in using a climate model for
116 long term projections, i.e., that the model represents the physics well enough to
117 trust the projections as long as the initial conditions are captured well. Ideally,
118 we would like to find a setting of the model parameters at which, when the model
119 is run without assimilation, the output most closely approximates the physical
120 behaviour and evolution of the climate system. At such a parameter choice, we
121 may then assimilate key data assuming that once assimilation is complete, the
122 free running model will “drift” back to its attractor slower than otherwise (as its
123 attractor is closer to the data than at any other parameter setting). In theory
124 then, with better parameter choices, short to medium term forecasts based on
125 data assimilation will be more accurate.

126 In this paper we present a statistical approach to climate model tuning using
127 an existing technique called history matching (Craig et al., 1996). It was initially
128 presented as a methodology for finding parameter settings of computationally ex-
129 pensive oil well models that led to output consistent with the observed history
130 of a well. It is currently used as a tool for quantifying parametric uncertainty
131 in computer models and has been applied to ice sheet models (McNeall et al.,
132 2013), intermediate complexity climate models Edwards et al. (2011) and to GCMs
133 (Williamson et al., 2013). The idea is for all parameters to be varied simultaneously
134 in the generation of a perturbed physics ensemble (PPE). The PPE is then used
135 to train emulators (fast statistical approximations to the climate model that give
136 a prediction of the climate model output for any setting of the parameters with
137 an associated uncertainty on the prediction) that are then used, in tandem with
138 observations, to cut out regions of parameter space that lead to models deemed
139 “too far” from the observations according to a robust geometric measure (we note
140 briefly here that if the model is extremely fast and uses no computational resource,
141 that emulators will not be required to map and reduce parameter space (Gladstone
142 et al., 2012). However, we present history matching as a method requiring emula-
143 tion to assist climate model tuning as these models have the opposite property).
144 Emulators have also been used in tuning exercises by Bellprat et al. (2012).

145 We argue that history matching is an effective and intuitive tool for tuning
146 and that it can be used to determine whether a perceived structural error actually
147 exists or if it can be corrected by changing the model parameters. History matching
148 can be applied on the most computationally expensive climate models and with
149 small ensembles, and represents a far simpler undertaking than a data assimilation
150 based approach. We illustrate the method by investigating a number of known
151 ocean circulation biases in HadCM3.

152 In particular we investigate the perceived structural bias in the Antarctic Cir-
153 cumpolar Current (ACC) strength. This current is known to be too strong in the
154 Hadley Centre climate models (Russell et al., 2006; Meijers et al., 2012); however,
155 we show that this may not represent a structural error at all. We show that by
156 jointly varying both ocean and atmosphere parameters together, it is possible to
157 find models that cannot be ruled out as having physical global surface air temper-
158 ature and precipitation using the metrics defined by Williamson et al. (2013), that
159 also have no ACC bias. We explore the properties of the ocean and atmosphere
160 circulations in one of these models and compare them to the standard HadCM3
161 and to observations. We use history matching to identify a region of parameter
162 space containing not implausible ACC strengths and further refine this region us-
163 ing a constraint on North Atlantic sea surface temperatures (SSTs). We investigate
164 qualitative features of the parameter space not ruled out by the observations of
165 these metrics and illustrate why ocean and atmosphere parameters must be varied
166 jointly in the coupled model when tuning.

167 In section 2 we briefly describe emulation and history matching and discuss its
168 implementation for tuning expensive climate models. In section 3 we use history
169 matching to search a subset of the HadCM3 parameter space with not implausible
170 SAT and precipitation profiles found by Williamson et al. (2013) for models with
171 not implausible ACC strength. We identify a subset of this space predicted to
172 contain not implausible ACC strengths and find some models therein. We inves-
173 tigate properties of the ocean circulation for a run without the usual ACC bias
174 and compare them to the standard HadCM3. In section 4 we include SSTs in the

175 sub tropical gyre into our history match and discuss features of the parameter
 176 space that has not been cut out. Section 5 contains discussion and the appendices
 177 present details of the emulation, the parameters varied in the ensemble and present
 178 further pictures.

179 2 Emulation and history matching

180 We write the climate model as the vector valued function $f(x)$ where x corresponds
 181 to a vector of climate model parameters. History matching requires an emulator
 182 for $f(x)$ to be fitted so that, for any setting of the parameters x , an expectation
 183 and variance for those elements of $f(x)$ we intend to compare with observations
 184 ($E[f(x)]$ and $\text{Var}[f(x)]$) may be computed from the emulator. There is a vast and
 185 growing literature on using ensembles to fit statistical emulators, so we don't go
 186 into mathematical details here. We refer the reader to Craig et al. (2001); Rougier
 187 (2008); Haylock and O'Hagan (1996); Sacks et al. (1989) and the book by Santner
 188 et al. (2003), for general information on building emulators; and to Rougier et
 189 al. (2009); Challenor et al. (2009); Sexton et al. (2011); Williamson et al. (2012);
 190 Schmittner et al. (2011); Lee et al. (2011); Williamson et al. (2013) and Williamson
 191 and Blaker (2014) for application of emulators to climate models.

192 Once an emulator is fitted so that we can compute $E[f(x)]$ and $\text{Var}[f(x)]$ for
 193 any x , history matching proceeds by ruling out choices of x as being inconsistent
 194 with chosen observational constraints, z , using an implausibility function $\mathcal{I}(x)$. A
 195 common choice is $\mathcal{I}(x) = \max_i \{\mathcal{I}_i(x)\}$ and

$$\mathcal{I}_i(x) = \frac{|z_i - E[f_i(x)]|}{\sqrt{\text{Var}[z_i - E[f_i(x)]]}}, \quad (1)$$

196 but others do exist (Craig et al., 1996; Vernon et al., 2010). Large values of $\mathcal{I}(x_0)$
 197 at any x_0 imply that, relative to our uncertainty, the predicted output of the
 198 climate model at x_0 is very far from where we would expect it to be if $f(x_0)$ were
 199 consistent with z . A threshold a is chosen so that any value of $\mathcal{I}(x_0) > a$ is deemed
 200 implausible. The remaining parameter space, $\{x \in \mathcal{X} : \mathcal{I}(x) \leq a\}$ is termed Not
 201 Ruled Out Yet (NROY). The value of a is often taken to be 3 following the 3
 202 sigma rule (Pukelsheim, 1994), which states that for any unimodal continuous
 203 probability distribution, at least 95% of the probability mass is within 3 standard
 204 deviations of the mean.

205 The form of $\text{Var}[z_i - E[f_i(x)]]$ will depend on any statistical model used to
 206 establish a relationship between observations of climate and output of the climate
 207 model. The most popular model, termed the 'best input approach' (Kennedy and
 208 O'Hagan, 2001) expresses the observations via

$$z = y + e$$

209 where y represents the underlying aspects of climate being observed and e rep-
 210 represents uncorrelated error on these observations (perhaps comprising instrument
 211 error and any error in deriving the data products making up z). The best input
 212 approach then assumes that there exists a 'best input' x^* so that

$$y = f(x^*) + \eta$$

213 where η is the model discrepancy (or structural error) and is assumed independent
 214 from x^* and from $f(x)$ at any x . Model discrepancy, being independent from any
 215 evaluation of the climate model, represents the extent to which the climate model
 216 fails to represent actual climate owing to missing or poorly understood physics,
 217 parameterisation schemes and the resolution of numerical solvers.

218 The best input approach has been used in studies with climate models by
 219 Murphy et al. (2009) and Sexton et al. (2011) and is described by Rougier (2007).
 220 The statistical model leads to

$$\text{Var}[z_i - E[f_i(x)]] = \text{Var}[e] + \text{Var}[\eta] + \text{Var}[f_i(x)]$$

221 where $\text{Var}[e]$ is the variance of the observation error, $\text{Var}[\eta]$ is model discrepancy
 222 variance and $\text{Var}[f_i(x)]$ is a component of the emulator for $f(x)$.

223 2.1 A tuning procedure: history matching in waves

224 History matching represents a formal statistical procedure for tuning climate mod-
 225 els by iteratively ruling out implausible regions of parameter space. We advocate
 226 tuning a climate model through a series of “waves” of history matching, where
 227 a “wave” involves running a new PPE in the current NROY space, building new
 228 emulators for each of the currently considered metrics and for a series of new
 229 metrics to be introduced for this wave, and using these emulators to further cut
 230 down NROY space. Structural errors are identified when a particular metric, once
 231 introduced, can rule out the whole space, indicating that, given all of the other
 232 metrics are NROY, the chosen metric cannot be reproduced to within the model
 233 discrepancy.

234 This approach has been demonstrated to be successful in other fields. For ex-
 235 ample, Vernon et al. (2010) demonstrate this procedure through five waves on
 236 a computer model simulating the evolution of galaxies after the big bang. After
 237 5 waves (with each ensemble containing 1000 different parameter settings) they
 238 found hundreds of computer model runs that were consistent with their observa-
 239 tions when, prior to the study it was thought that no such parameter settings
 240 existed. Given this success in other fields, we believe that it is highly likely that
 241 at least some of the perceived “structural errors” in modern GCMs will be elimi-
 242 nated by history matching without compromising model performance with respect
 243 to other physically important metrics.

244 The crucial decision to be made by the modellers when using history matching
 245 in this way, is what metrics should be used to tune the model and in what order
 246 should they be applied. There are aspects of real world physics that we know, a
 247 priori, that the model does not capture. For example, sub grid scale processes such
 248 as eddies in HadCM3, or convective plumes in a $1/4^\circ$ model. These are definitely
 249 not part of the model, and so if we were to history match to them, we would rule
 250 out the whole parameter space. Hence the choice of metric is important.

251 Further, the order in which they are introduced through the different waves is
 252 also important. Note that within a given wave, all metrics have the same level of
 253 importance. If a parameter choice is ruled out because of one metric, it is ruled
 254 out no matter if it is NROY with respect to others in the same wave or not.
 255 However, the wave at which each metric is introduced should reflect the order of
 256 importance of any particular metric when it comes to trusting the output of a

257 climate model. For example, it may be that the model must have a reasonable
258 global SAT profile and that this is more important than its AMOC strength. In
259 this case, by first matching to SAT in wave 1, then searching the space of models
260 with NROY SAT profiles for reasonable AMOC strengths in wave 2, we only search
261 a sub-space of models for good AMOC strengths. If we fail to find any, we would
262 declare that AMOC was a structural error in the model. However, it might be that
263 certain parameter choices that lead to poor SAT profiles do have not implausible
264 AMOCs. By choosing the order in which metrics are introduced over successive
265 waves, we effectively define what it means for the model to have structural error.
266 Models are currently tuned by comparing various metrics to observations by the
267 modellers. Hence there is already an implicit sense of what metrics are important
268 and which are more important than others.

269 An important issue if adopting our approach to tuning through history match-
270 ing is that of stopping rules. How many wave of refocussing and history matching
271 are required before one or more of the NROY models can be adopted as a tuned
272 run (or collection of them)? The answer to this will be very problem dependent,
273 however we can offer guidelines. If the set of chosen metrics is fixed then once the
274 emulator variance is a great deal smaller than the denominator in the implausibil-
275 ity calculation, then it is unlikely that further waves will change the implausibility
276 very much, and further ensembles for history matching purposes may be consid-
277 ered to be an inefficient use of resources. In these situations, parameter choices
278 will generally have low implausibility because we are relatively sure that the model
279 when run at those choices is genuinely close to the observations (with respect to
280 the uncertainty and model discrepancy). How many waves are required in order to
281 reach this situation will depend on a number of factors. These include, the avail-
282 able ensemble size at each wave, the complexity of the metric and its behaviour in
283 parameter space (is it easy to emulate with few runs, or is a lot of data required
284 to capture the parameter dependencies?), the size of the sub-volume of parameter
285 space that, if we did run the model there, would be close enough to the observa-
286 tions (which depends on their uncertainty on on the model discrepancy) and the
287 wave number at which new metrics are introduced.

288 Before we move on, we address the issue of specifying model discrepancy. A
289 reader might object that we are advocating a methodology for locating structural
290 errors that requires us to know already what the structural errors are by providing
291 a model discrepancy variance. Certainly, if model discrepancy variance for any
292 metric can be specified or estimated by experts, then history matching can proceed
293 straightforwardly. However, when tuning we do not expect this to be the case. If it
294 is not, we can treat the discrepancy variance as our tolerance to structural error.
295 This enables us to explore parameter space and discover whether or not regions
296 containing parameter settings that are not inconsistent with the observations we
297 would like to match to exist with respect to different tolerances to this error.

298 The notion of specifying a tolerance to error should not be unfamiliar to model
299 developers tuning their climate models, where the goal is often to tune components
300 of the climate model so that they are “close to” observations. How close is accept-
301 able will be known to the modellers who are often varying one or a handful of
302 parameters thought relevant to that process at any one time until an “acceptable”
303 setting of the model parameters is found (or it is thought that a structural error
304 exists). Hence, part of the definition of a structural error, is what the tolerance to
305 model error is. For example, we might be more tolerant to errors in the location

306 of the Gulf Stream in a 2° model than we would be to errors in its global mean
307 temperature.

308 2.2 Computationally expensive models

309 One potential objection to adopting a rigorous statistical approach to climate
310 model tuning that uses PPEs is that the latest climate models are too expensive
311 to run, so that PPEs large enough to build emulators with cannot be obtained.
312 This is not a problem for history matching.

313 History matching only requires an emulator for a climate model. Though one
314 effective way to emulate a climate model is to use a large PPE, it is not the only
315 way. In fact an emulator can be built for a model for which you have no data
316 at all (Goldstein and Rougier, 2009; Williamson and Goldstein, 2013)! The most
317 practically effective way to build an emulator for a slow, expensive climate model is
318 to use a large PPE on a coarse resolution version of it. For example, as mentioned
319 earlier, the ocean component of HadGEM3 is the 0.25° resolution version of the
320 NEMO ocean model. This model can also be run much more quickly at 2° and 1°
321 resolution.

322 The idea is to use a large ensemble of coarse resolution models and to write
323 down an emulator for the expensive model as a function of the emulator for the
324 coarse version. Note that this is an emulator and that we need no runs of the
325 expensive model to construct it. Though this emulator is likely to have large un-
326 certainties on the predictions it makes, particularly when changes in resolution lead
327 to changes of parameterization schemes, it can then be efficiently tuned using very
328 small ensembles from the expensive model in order to reduce these uncertainties.

329 To be clear, we are not suggesting that a 1° model run at some parameter
330 choice x_1 would be informative for the 0.25° output at x_1 . For example, some
331 parameters in x_1 may not be needed in the 0.25 degree model, and other new
332 parameters may be required to run it. Our claim is that because the models are all
333 simulating climate, then unless we claim that one of the versions has no skill even in
334 reproducing large scale features of the climate (such as global mean temperatures
335 or circulations), the different parameter spaces must be related. Large ensembles
336 of the coarse model and small ensembles of the high resolution model can be
337 used to establish this relationship statistically through an emulator for the high
338 resolution version. New coarse experiments can reduce some of the uncertainty in
339 this emulator as can new high resolution experiments, and the emulator can be
340 used in history matching as with any other emulator.

341 For example, Williamson et al. (2012) emulate 200 year time series of the
342 Atlantic Meridional Overturning Circulation (AMOC) in the coupled HadCM3
343 using a PPE with just 16 members and a large ensemble of a coarse version called
344 FAMOUS. Given the prior emulator for the expensive model, one can use its un-
345 certainty specification to aid experimental design decisions so that a finite budget
346 of model evaluations can be spent on removing as much parameter space as pos-
347 sible. This will be more efficient than the “one factor at a time” type of approach
348 that is currently used. For more information on emulating expensive models us-
349 ing coarse resolution versions see Cumming and Goldstein (2009); Kennedy and
350 O’Hagan (2000); Le Gratiet (2014) and Williamson (2010).

Time and budget constraints prevent us from obtaining further ensembles of HadCM3 for this study, so that we cannot demonstrate iterative tuning for this model. However, in the rest of the paper we demonstrate the potential effectiveness of history matching for tuning climate models by using further constraints on the NROY space found in Williamson et al. (2013). These constraints are designed both to remove regions of parameter space with poor ocean circulations (including the standard HadCM3) and regions with the observed SST biases. Following this second history match, we can plot 1 and 2D projections of NROY parameter space and find which parameters drive the majority of the reduction of parameter space.

3 The Antarctic Circumpolar Current in HadCM3

At a recent workshop on ocean model discrepancy, a group of oceanographers and statisticians discussed key processes in the ocean that drive the AMOC and that would have to be modelled correctly in order for them to have confidence in the modelled transient response of the AMOC to CO₂ forcing. Of the processes mentioned, some of those deemed more important included location and strength of the sub-polar and sub-tropical gyres, temperature and salinity in the sinking regions in the North Atlantic and the strength of currents in the Southern Ocean. These discussions also led to a number of ocean processes in HadCM3 that were thought to impact upon AMOC strength being identified as having “known structural biases”. We are motivated in this illustration of history matching as a tool for tuning climate models by investigating the nature of the biases that our experts deemed influential on the AMOC. We begin with the ACC strength.

ACC strength (Sv, $1 \text{ Sverdrup} = 1 \times 10^6 \text{ m}^3 \text{ s}^{-1}$), measured across Drake Passage, is an ocean transport that has proved difficult to capture accurately in AOGCMs. In the multi-model ensemble used to support the Intergovernmental Panel on Climate Change’s fourth assessment report (IPCC-AR4 Solomon et al., 2007) known as CMIP3 (Coupled Model Intercomparison Project phase 3 Meehl et al., 2007), the range of ACC transports given by the then state of the art climate models (including HadCM3) was huge compared to the observations and their associated uncertainty (134 ± 15 to 27 Sv though this error is misquoted as being 11.2 Sv Cunningham et al., 2003). The CMIP3 models ranged from -6 to 336 Sv, but, perhaps more surprisingly, only two of the models returned an ACC strength consistent with the observations (see Russell et al., 2006, for details). The CMIP5 models (Meijers et al., 2012) fare a little better with a range of $90 - 245$ Sv but still only 2 models consistent with the observations. The Hadley centre models are all too strong in CMIP5.

If an overly strong ACC strength represents a structural error in HadCM3, this would imply that it is not possible for HadCM3 to simulate a realistic climate with an ACC strength close to observations. We investigate this possibility using a large PPE of HadCM3 runs described below.

3.1 The ensemble

We designed a large PPE on the coupled, non-flux-adjusted, climate model HadCM3 (Gordon et al., 2000; Pope et al., 2000). This ensemble varied 27 parameters con-

394 trolling both the model physics in the atmosphere and ocean of HadCM3 and was
395 generated using Climate Prediction Dot Net (CPDN, <http://climateprediction.net>). CPDN is a distributed computing project through which different climate
396 models are distributed to run on personal computers volunteered by members of
397 the public. A copy of the model, along with a specific prescribed setting of the
398 model parameters, is downloaded by the “client” computer, where it runs in the
399 background using any spare computing resources available. Data is returned to
400 CPDN where it is stored and made available for access by the general public.
401

402 The ensemble consists of a 10,000 member design in the chosen parameters
403 submitted in April 2011. At the time of writing there are over 3500 unique ensemble
404 members that have completed 120 years of integration with preindustrial boundary
405 conditions. Information on the design of the ensemble can be found in Williamson
406 et al. (2013) and Yamazaki et al. (2012). A comprehensive list of the parameters
407 varied appears in appendix B.

408 Though we were able to use CPDN to generate a very large ensemble, our
409 method does not rely on ensembles that are so large that only a program similar to
410 CPDN would render it practical. Loeppky et al. (2009) suggest that a good rule of
411 thumb for ensemble size, in the absence of information from previous experiments,
412 is 10 times the number of parameters. Sexton et al. (2011) built GCM emulators
413 for UKCP09 using slightly larger ensembles than this using the UK Met Office
414 supercomputer. As discussed in section 2.2, emulators, and thus history matching,
415 can be achieved using extremely small ensembles of very expensive models, if large
416 ensembles (e.g. large enough to satisfy the rule of thumb suggested by Loeppky et
417 al. (2009)) on related coarser resolution models are available.

418 3.2 NROY space

419 Williamson et al. (2013) perform a history match on HadCM3 using 4 observational
420 metrics to cut out over half of the original parameter space. The NROY space for
421 HadCM3 derived in Williamson et al. (2013) consists of all those parameter settings
422 that couldn’t be ruled out using global mean surface air temperature (SAT), global
423 mean precipitation (PRECIP), the global mean surface air temperature gradient
424 (SGRAD) and the global mean seasonal cycle in surface air temperature (SCYC).
425 Hence any parameter setting in NROY space already has a not implausible global
426 mean surface air temperature profile and global mean precipitation with respect
427 to the chosen constraints.

428 3.3 NROY ACC

429 In order to further constrain NROY space using the ACC strength by history
430 matching, we require an emulator for the ACC strength as well as an observational
431 error variance and a discrepancy variance. We describe the emulator for ACC
432 strength in appendix A and we interpret the Cunningham et al. (2003) range ($134 \pm$
433 15 Sv) as 3 standard deviations using the 3 sigma rule (Pukelsheim, 1994). This
434 gives $\text{Var}[e] = 25$. We note here that there are many ways one might interpret the
435 error quoted in data range statements such as this. One is that the range represents
436 hard boundaries on the value of the true process. Under this interpretation our

437 representation of the range as 3 standard deviations leads to a larger error variance
438 than necessary and so less parameter space ruled out through history matching.
439 The interval might be viewed as a confidence interval for the true value of the data,
440 and our interpretation is consistent with the quoted range as a 95% confidence
441 interval under the assumption that the underlying distribution of the observations
442 is unimodal (Pukelsheim, 1994). A third way might be to interpret the quoted
443 range as 1 standard deviation, however to use this interpretation here would imply
444 that the search for models with an ACC strength any closer than 45 Sv (3 standard
445 deviations) from the observations would be overfitting (even with a perfect model
446 and perfect emulator). We know from conversations with NEMO developers and
447 from the discussion of the performance of the CMIP5 models that the data are
448 treated as being more accurate than this and that the field looks for models that are
449 within the quoted data range. Treating the quoted range as 3 standard deviations
450 is consistent with the desire to search for model runs that meet this constraint.

451 We specify zero tolerance to climate model error via a model discrepancy vari-
452 ance of 0, so that we demand that the model output lies within the range of the
453 observation uncertainty. This assumption is not one we would make if our goal
454 were to tune HadCM3, as we do have tolerance to model error. This study aims to
455 explore the capabilities of HadCM3, hence, in specifying zero tolerance to model
456 error, we are testing to see if the model is capable of replicating the observations
457 whilst having a reasonable global temperature and precipitation profile. If we were
458 actually looking to tune the climate model these tolerances would not be zero and
459 may well be correlated across constraints as a modeller may tolerate more error
460 in one type of constraint (e.g. AMOC strength) in favour of less error in another
461 (e.g. SST). One argument, for example, for including a non-zero model discrep-
462 ancancy variance here, is that our simulations are preindustrial and hence we might
463 want to account for climate change in the observations. However, we might also
464 view this uncertainty as part of the observation error.

465 From Williamson et al. (2013) we know that 56% of the defined parameter space
466 is removed using the first 4 constraints. Demanding not implausible ACC strength
467 reduces the remaining space by 90.4% leaving just 4.3% of the parameter space
468 not ruled out yet. We explore the properties of the parameters in this NROY space
469 in section 4, however, we note that we have ensemble members that satisfy each of
470 our 5 constraints and focus the rest of this section on exploring the behaviour of
471 one of these models in particular. Figure 1 plots the mean ACC strength for the
472 final decade of every member of our ensemble against the mean AMOC strength.
473 Dashed lines represent the lower and upper bounds on the observations. Points
474 outside of this box have ACC strength, AMOC, or both outside of the range given
475 by the observations and may be thought of as having unphysical ocean circulations.
476 We colour points ruled out by our wave 1 history match in Williamson et al. (2013)
477 in grey and add the NROY members from this analysis in cyan. We colour those
478 points NROY to the additional constraint of ACC strength in dark blue. Standard
479 HadCM3 is plotted as the red triangle.

480 From this plot we can see that we have a number of not ruled out yet ensemble
481 members with a not implausible ACC strength. We also see that the standard
482 HadCM3 (plotted as the pink triangle) has an overly strong ACC. Note that
483 many (more than the 5% expected) of the now NROY ensemble members (the
484 blue points) still have ACC strengths outside of the observation range. This is a
485 feature of history matching with emulators. We only rule parameter choices out

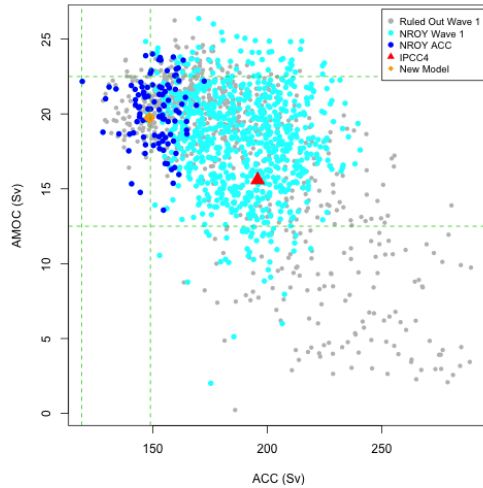


Fig. 1 The ACC (Sv) through Drake Passage plotted against the AMOC (Sv) at 26°N in our HadCM3 ensemble. The dashed lines represent upper and lower bounds on observations of ACC and AMOC. We colour those ensemble members ruled out by the initial history matching in Williamson et al. (2013) grey, with NROY members from this wave in cyan. Members in NROY space when adding ACC strength as a constraint appear in blue. The standard HadCM3 is highlighted on this plot as the red triangle. The new run we examine in further detail in this section is shown as the orange diamond.

486 if we are sure they lead to unphysical circulations, and part of our uncertainty
 487 comes from the quality of our emulator (and the predictability of the model). Our
 488 particular choice of emulator (see appendix A) does not interpolate the ensemble
 489 members and give zero variance at those parameter locations because we have
 490 accounted for internal variability in our modelling. So we expect to see parameter
 491 choices, and therefore existing ensemble members, that are predicted to be (or
 492 have been observed to be) outside the target data range as NROY because our
 493 uncertainty (driven by internal variability in this case) is such that we can't be sure
 494 that they really do lie outside the data range for any setting of initial conditions.

495 3.4 Comparison with the standard HadCM3

496 The NROY members identified with the ACC constraint all pass the large scale
 497 constraints imposed by the initial history match (SAT, PRECIP, SGRAD, SCYC).
 498 We now examine the state of the climate of one of these members (the orange
 499 diamond in figure 1) in more detail. The ensemble member selected for more
 500 detailed analysis is typical of the NROY members identified in blue on figure 1,
 501 and whilst for any individual metric it is possible to select other equally good or
 502 better example (we present one such example for the barotropic streamfunction
 503 in appendix C), this member represents a good compromise of all the metrics we
 504 examine.

505 Although there are no direct observations of the barotropic streamfunction,
 506 comparison with other models and reanalyses, as well as observations of sea surface

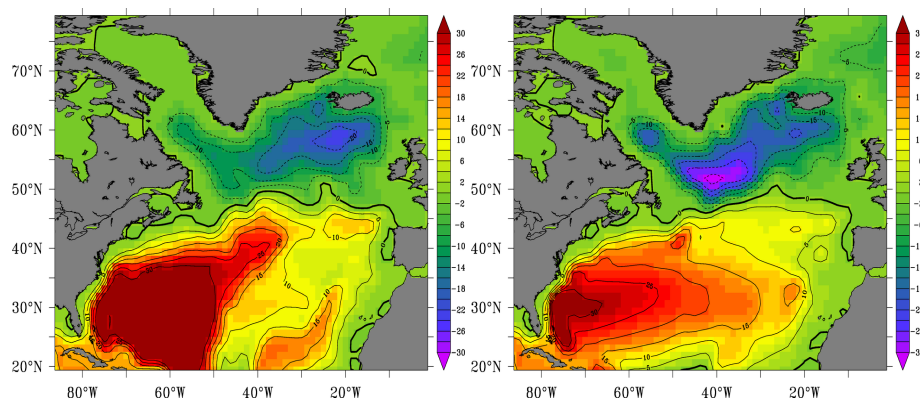


Fig. 2 The barotropic streamfunction (BSF) for the standard HadCM3 (left) and an ensemble member with realistic ACC strength (right). Units are Sv. The domain is smoothed spatially to remove grid point noise.

507 height from altimetry can be made. These comparisons indicate that the subpolar
 508 gyre is too far east in the standard HadCM3 (see figure 2), with its centre located
 509 around 25°W , 55°N (south of Iceland), and the subtropical gyre is too broad
 510 and diffuse, with most of the southward return flow occurring in a narrow band
 511 between $50\text{--}55^{\circ}\text{W}$. In comparison, the alternative model with the more realistic
 512 ACC has a subtropical gyre which is much more tightly constrained to the western
 513 boundary and has more uniform southward return distributed across the rest of
 514 the basin. It also has a westward shifted subpolar gyre with a maximum transport
 515 located around 45°W , 55°N (south of Cape Farewell), though there is still a strong
 516 gyre near Iceland. Some of the ensemble members do exhibit weaker subpolar gyre
 517 circulation south of Iceland (see appendix C), although we note that modifications
 518 to the bathymetry in HadCM3 (excavation of the sills and a submerged Iceland)
 519 may force the model to produce an unrealistic eastward extension of the SPG
 520 circulation.

521 We continue our examination of the chosen NROY model with an assessment of
 522 the SST, SSS and circulation represented in the North Atlantic and Nordic Seas,
 523 a region considered very important because of the formation of North Atlantic
 524 Deep Water (NADW), which forms the lower limb of the AMOC. Figure 3 plots
 525 sea surface salinity (SSS) (left panels) and sea surface temperature (SST, right
 526 panels) anomalies with respect to the observations (Ingleby and Huddleston, 2007)
 527 for the standard HadCM3 (top panels) and the chosen NROY model in the bottom
 528 panels. A number of supposed “structural errors” in HadCM3 can be identified on
 529 the upper panels and have been improved by the alternative parameter choice. For
 530 example, the standard HadCM3 has a fresh bias of -0.5 PSU to -1.5 PSU in the
 531 subpolar gyre and along the Greenland coast. Extending out into the Norwegian
 532 sea the fresh bias can exceed -3 . The fresh bias in the subpolar gyre is not present
 533 in the alternative model and the one along Greenland is halved. We also improve
 534 the salty bias which extends all the way down the eastern boundary and across the
 535 Atlantic at $20\text{--}25^{\circ}\text{N}$. These improvements are at the expense of a slight freshening

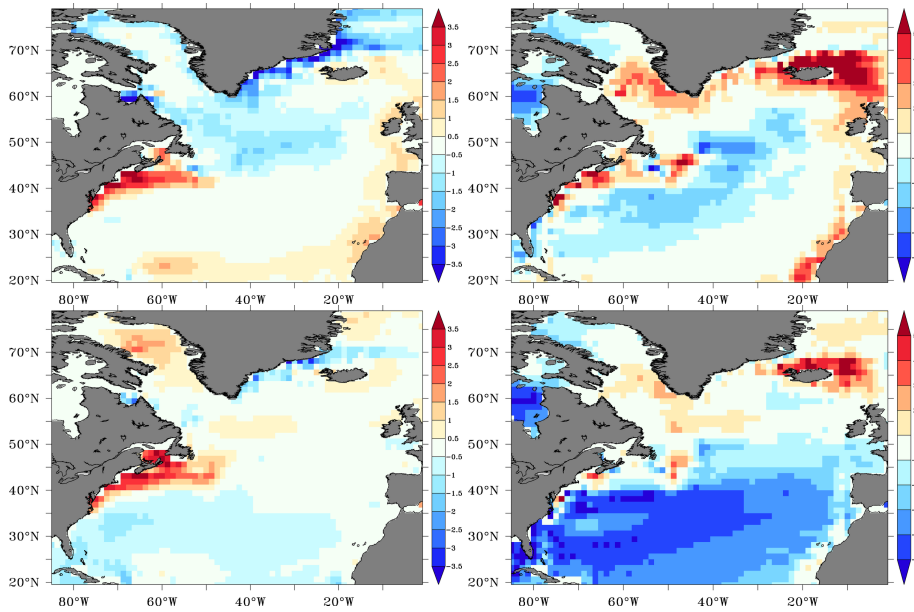


Fig. 3 Sea surface salinity (SSS) anomalies (PSU, left panels) and sea surface temperature (SST) anomalies ($^{\circ}\text{C}$, right panels) for the standard HadCM3 (top 2 panels) and the ensemble member with realistic ACC strength (bottom two panels). Both sets of anomalies are calculated as the difference of the mean of the last ten years in our ensemble and EN3 (Ingleby and Huddleston, 2007).

536 of the subtropical gyre. However, it is arguably more important that the salinity
 537 is correct in the AMOC sinking regions.

538 Note that both models exhibit the same positive salinity anomaly at the region
 539 of the Gulf Stream separation, indicative of a structural error which arises because
 540 of the model resolution. We believe that this anomaly is not possible to correct
 541 in HadCM3 by tuning parameters. The SST is also closer to observations in the
 542 North Atlantic sinking regions. Most notably the warm bias around Iceland, and
 543 extending west round Greenland and into the Labrador sea is reduced. However,
 544 these improvements are accompanied by the development of a larger and stronger
 545 cold bias in the sub tropical gyre. This cold bias is undesirable, and we would try
 546 to address this in the next wave. Many, though not all, of our NROY members
 547 have a similar cold bias in the sub tropical gyre. We believe that this bias (and
 548 the others introduced, see below) can be tuned out as its strength and sign vary
 549 across our ensemble, unlike, say, the salinity bias in the gulf stream, which was
 550 present in every ensemble member.

551 We can examine the circulation of the subtropical gyre in more detail. Fig-
 552 ure 4 shows a cross section of the meridional velocity at 26°N for the standard
 553 HadCM3 (top) and the alternative model with improved ACC (bottom). There
 554 is no available climatology to compare with for this field, but we refer the reader
 555 to observations presented in Johns et al. (2008) (their Fig. 7), and to the high
 556 resolution model climatology validated and used by Meinen and Garzoli (2014)
 557 (their Fig. 6). In the standard model the deep western boundary current is too

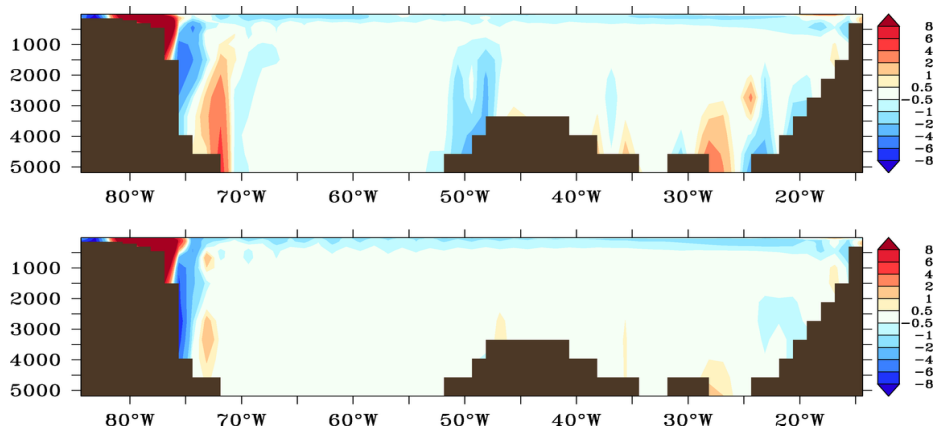


Fig. 4 Cross section of the meridional velocity (ms^{-1}) at $26^{\circ}N$ for the standard HadCM3 (top) and the ensemble member with realistic ACC strength (bottom). Red indicates northward flow, and blue indicates southward flow. Units are cm/s .

558 broad and shallow at the western boundary and there is a substantial northward
 559 transport below 2500 m between $72-74^{\circ}W$. The intense southward transport indicated by the tight contours in figure 2 ($50-55^{\circ}W$) can be identified as a strong
 560 return flow at the western flank of the mid-Atlantic ridge and there are additional
 561 spurious transports on the eastern boundary. The alternative model with improved
 562 ACC transport exhibits a more physical deep western boundary current, which is
 563 tightly constrained to the western boundary. It does not have the large return
 564 flow on the western flank of the mid-Atlantic ridge and transports in the eastern
 565 basin are more realistic. By simply finding a NROY model with a more physical
 566 ACC transport through Drake Passage, we have found a model with an improved
 567 representation of the ocean circulation in the North Atlantic. However, we must
 568 also verify that these improvements have not arisen at the expense of the model
 569 developing serious problems elsewhere in the global climate.
 570

571 Figure 5 compares the global SSS anomaly (left panels) and SST anomaly
 572 (right panels) fields for the standard HadCM3 (top panels) and the improved
 573 ACC member (bottom panels). SSS is improved almost everywhere outside of the
 574 Arctic Ocean. We note that HadGEM1, the successor of HadCM3, showed similar
 575 improvements to SSS globally, also replacing the Arctic SSS dipole anomaly with
 576 a pan Arctic positive anomaly (Johns et al., 2006). The SSS anomalies were also
 577 improved in CHIME (Megann et al., 2010), a coupled model closely related to
 578 HadCM3 where the ocean component was replaced by HYCOM and interestingly
 579 CHIME also exhibits the same pan Arctic positive SSS anomaly.

580 SST anomalies still present a problem. The alternative model shows the same
 581 gradient in SST anomalies, with the northern hemisphere exhibiting a cold bias
 582 and the southern hemisphere exhibiting a warm bias. The North Pacific cold bias is
 583 much stronger in the NROY simulation, exceeding $5^{\circ}C$, and a cold bias associated
 584 with excessive equatorial upwelling can be seen in the Pacific. South of $20^{\circ}S$ the
 585 alternative model shows substantial improvement, with the warm biases being
 586 reduced both in area and amplitude. Interestingly, the North Pacific cold bias is

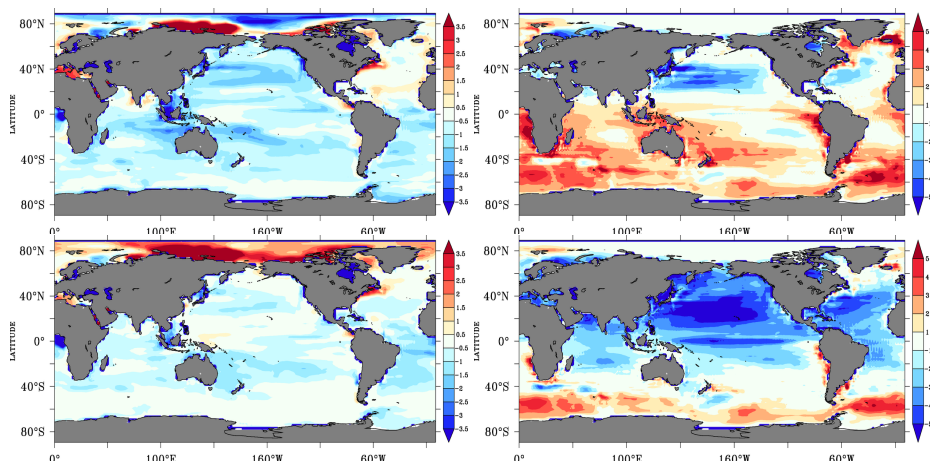


Fig. 5 Left panels are global SSS anomalies (PSU) and right panels are global SST anomalies ($^{\circ}\text{C}$), both with respect to EN3 (Ingleby and Huddleston, 2007) climatology. The top panels are for the standard HadCM3 and the bottom panels are from an ensemble member with realistic ACC strength.

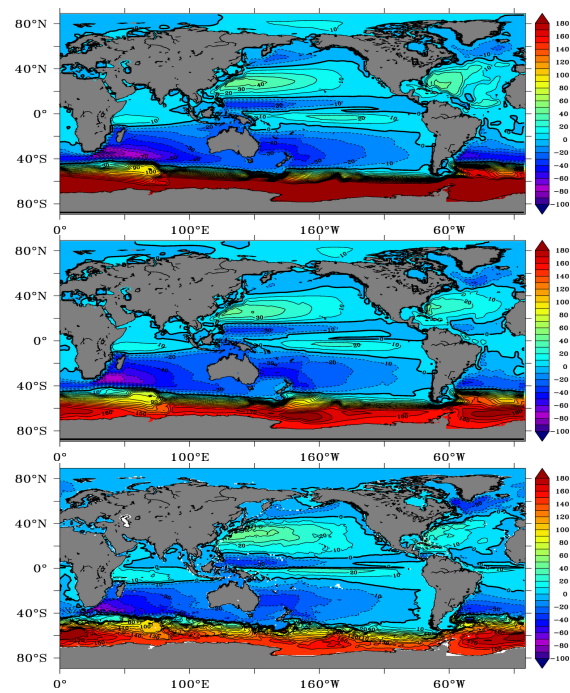


Fig. 6 Global barotropic streamfunction for the standard HadCM3 (top), the improved ACC run (middle), and for a NEMO ORCA025 ($1/4^{\circ}$) ocean only simulation

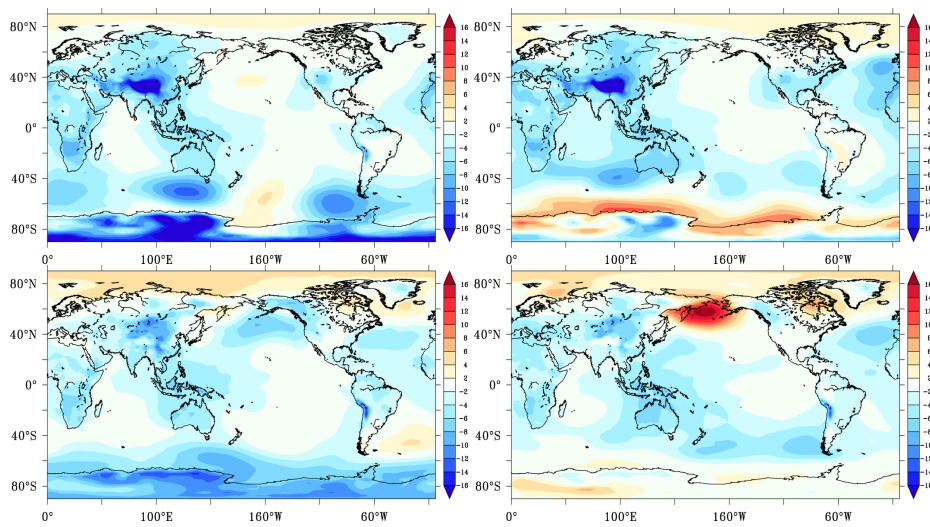


Fig. 7 Sea level pressure (SLP) anomalies ($mbar$) with respect to ERA-40 (Uppala et al., 2005) 1960-1990 climatology for summer (left) and winter (right). Anomalies for the standard HadCM3 are shown in the top panels and the bottom panels are anomalies from the improved ACC strength run.

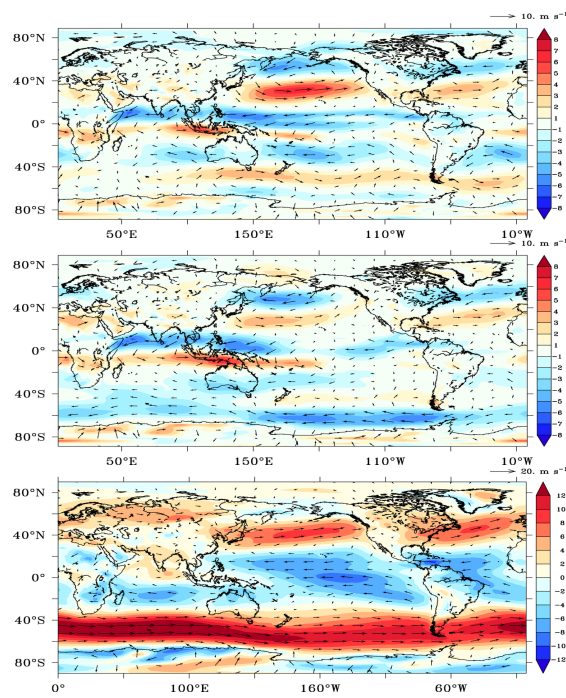


Fig. 8 Global wind field anomalies from the ERA-40 reanalysis at 850hpa (ms^{-1}) for the standard HadCM3 (top) and the improved ACC run (middle). Shading represents the zonal wind anomaly. The ERA-40 reanalysis climatological mean is also presented for reference (bottom)

587 present in both HadGEM1 (Johns et al., 2006) and HiGEM (Shaffrey et al., 2009)
588 is much less evident in CHIME (Megann et al., 2010), but the North Pacific is
589 still recognisably biased cold compared with the global mean. The gradient in
590 SST anomalies (generally cold in the northern hemisphere with warmer biases in
591 the southern hemisphere) may be indicative of different biases in the air-land and
592 air-sea interactions, suggesting that the northern and southern hemisphere biases
593 could be controlled by different parameters and that there is therefore scope to
594 reduce the slope and improve the bias overall.

595 HadCM3 is a rigid lid model, so there isn't a sea surface height field to com-
596 pare with observations, and the rigid lid pressure is not saved as a time mean
597 output. We do, however, have the time mean barotropic streamfunction, which we
598 can compare with state-of-the-art high resolution ocean only simulations forced
599 with observed surface fluxes. Comparing the global barotropic streamfunction in
600 the standard and improved ACC simulations to a simulation of NEMO ORCA025
601 ($1/4^\circ$) referred to as N206 (figure 6; see Blaker et al. (2014) for a detailed descrip-
602 tion of N206), substantial improvements exist over both the Atlantic and Southern
603 Oceans. N206 has a mean ACC transport of 139 Sv, and an AMOC transport of
604 14.9 Sv. The barotropic streamfunctions of the two HadCM3 simulations are sim-
605 ilar over the Indian and Pacific oceans, with the standard HadCM3 producing a
606 slightly stronger Pacific subtropical gyre and Kuroshio, as seen in N206.

607 There are improvements in the SLP (Figure 7), particularly over the South-
608 ern Ocean and Antarctic continent, but also over land across much of the globe,
609 most notably over the Himalayas. The improvements in SLP are closely linked
610 to the wind field (figure 8). Comparing the 850 hpa wind anomalies for the two
611 simulations (top two panels), differenced from the ERA-40 (Uppala et al., 2005)
612 reanalysis 1960-1990 mean (bottom panel) there are improvements in the mean
613 wind field over much of the globe, with the improved ACC model showing more
614 realistic wind strength over northern extent of the Southern Ocean, and better
615 easterly winds over the tropical Pacific and Atlantic. The zonal wind over the
616 southern extent of the Southern Ocean is anomalously weak in the improved ACC
617 run. Both simulations exhibit a too-zonal storm track over the North Atlantic,
618 whilst in the Pacific the storm track is too far south from around 150°E to 180°E .
619 The standard HadCM3 storm tracks are stronger than the climatology, whilst the
620 storm tracks simulated by the alternative model with improved ACC are slightly
621 weaker. The model with improved ACC representation displays a localised, strong
622 positive SLP anomaly in the North Pacific between Kamchatka and Alaska, related
623 to the weaker and more southerly Pacific storm track. This may be considered un-
624 desirable in a model used for UK weather and climate prediction if it affects the
625 characteristics of the northern hemisphere storm track, however the limited data
626 available from the CPDN ensemble means we cannot investigate this more closely.

627 Comparing precipitation anomalies between standard HadCM3 and the alter-
628 native model we see reductions in the error almost everywhere, particularly over
629 the maritime continent and along the ITCZ, but also in the subtropical regions
630 (figure 9). However, we note that there are large uncertainties in precipitation
631 climatologies so these improvements should be regarded with caution.

632 In figure 10 we plot the Meridional Heat Transport (MHT) and the AMOC
633 for every member of our ensemble, and highlight the standard HadCM3 (red) and
634 the improved ACC (blue) members. The grey lines are ruled out (RO) ensemble
635 members and the other colours correspond to the value of entcoef, the parameter

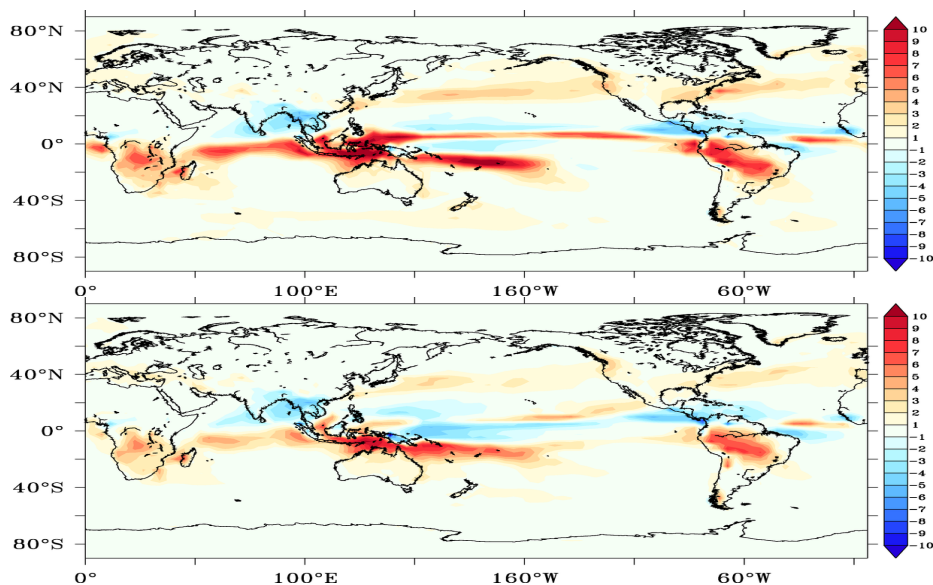


Fig. 9 Precipitation anomalies (*mm/day*) from ERA-40 (Uppala et al., 2005) 1960-1990 climatology for the standard HadCM3 (top) and the improved ACC model (bottom).

636 which governs the convective cloud entrainment rate. We see that both the MHT
 637 and AMOC are stronger and closer to observations. We also note that the vari-
 638 ability in the AMOC is larger for the improved ACC run, an observation that was
 639 true of each of the improved ACC runs we looked at.

640 4 Refining the search

641 Through jointly varying atmosphere and ocean parameters and using history
 642 matching we have found a region of parameter space with a predicted not implausi-
 643 ble global mean temperature profile, global mean precipitation and ACC transport
 644 through Drake Passage. This region of parameter space represents only 4.3% of
 645 the original space. Within this space we have found a version of HadCM3 that
 646 outperforms the standard version in many aspects of its ocean and atmosphere,
 647 and is arguably a better model. However, some might argue that the improvements
 648 to certain aspects of the ocean and atmospheric circulation have come at the too
 649 high price of exacerbating the cold bias in the north Pacific. To head off such ob-
 650 jections, we are not claiming to have found the “best” HadCM3 in any sense, nor
 651 that if the goal were to tune HadCM3, one should choose ACC transport through
 652 Drake Passage as a primary metric to cut down parameter space. It might be
 653 that, if a certain combination of metrics were used and ACC were left out, then
 654 something close to standard HadCM3 would be found (though our results make us
 655 doubt this). Instead, we claim that this method will find models that will exhibit
 656 improvements in key metrics chosen by the modeller if all of those metrics can be
 657 improved by changing the model parameters.

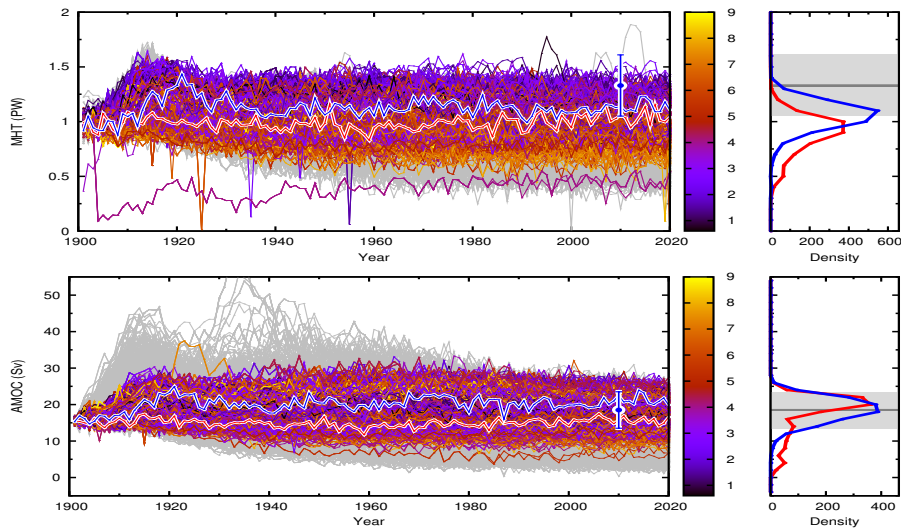


Fig. 10 The Meridional Heat Transport (top panel) and AMOC time series for each member of our ensemble. Grey runs have been ruled out by history matching. NROY runs are coloured by their value of entcoef. The highlighted red line in each image is the time series for the standard HadCM3 and the highlighted blue line is the time series for our alternative with realistic ACC transport. The curves on the right shows the unscaled density of the final year data for the NROY part of the ensemble (blue) and the ruled out part of the ensemble (red). The blue point represents the observations and includes error bars (also reproduced as grey shading in the unscaled density plots on the right hand side).

658 Further, the first model that is found within the NROY space that satisfies
 659 current metrics will not be the optimal version of the climate model unless we have
 660 been extremely fortunate to have hit the exact optimal parameter setting within
 661 our 27 dimensional space of continuous parameters. Instead, this model provides
 662 insight into which metrics have been improved and which must be used in further
 663 history matching (as long as this is physically appropriate). In our application to
 664 HadCM3, we would want to refine our model search to only include models that
 665 correct the ACC bias whilst simultaneously improving the representation of global
 666 SST.

667 In order to refine our search for HadCM3's with fewer, less serious structural
 668 biases using history matching, the process is to first select metrics on which to
 669 match, to emulate these within the current NROY space, and to use the results
 670 to design a new ensemble within the latest NROY space, to refocus our statistical
 671 models, then repeat in order to converge either on a set of models that repro-
 672 duces all specified metrics to within the chosen error tolerances, or upon a number
 673 of metrics that cannot be simultaneously reproduced and can thus be correctly
 674 identified as structural errors. Time and budget constraints have left us unable to
 675 run any further ensembles of HadCM3 with CPDN, however, we can perform the
 676 first part of this task, in order to gain insight into features of the NROY space of
 677 HadCM3 when we require that SAT, Precip, ACC and SST profiles are not un-
 678 physical. A future project might then seek to populate this space with models in

679 order to investigate their properties further, to further cut down NROY space, and
680 to look at transient simulations of the most physical looking models if appropriate.

681 To do this, we include the SST anomaly in the sub-tropical gyre as a chosen
682 metric. We could have included features of the Pacific SST or even certain spatial
683 patterns as metrics for history matching, however, it was felt that the most crucial
684 region to get right in order to have confidence in the AMOC, was the North
685 Atlantic, and our model exhibits a large SST bias there too. It was also felt that
686 correcting this bias, if possible, would simultaneously improve the temperature of
687 the North Pacific. We define our metric to be the mean SST in a box from 70°W to
688 30°W and from 26°N to 36°N . The “improved” HadCM3 we found has anomalies
689 up to 5°C in this region. We specified a tolerance to error of half of that, so that
690 our model discrepancy has a 3 standard deviation range of 2.5°C . Our discrepancy
691 variance is therefore 0.69. The region of the Atlantic we are assessing here is very
692 well observed compared with global SST, so the observation error variance is likely
693 to be low (calculations based on Ingleby and Huddleston (2007) estimated the
694 observation error variance on global mean SST observations over a 30 year period
695 as around 0.003). We therefore ignore observation uncertainty for this constraint,
696 taking the view that it is negligible relative to model discrepancy variance.

697 4.1 Results

698 In this section we explore the NROY space left when history matching to our 4
699 prior constraints as well as the ACC strength and the SST in the sub-tropical gyre.
700 To do this we evaluate implausibilities, via equation (1), for millions of untried
701 points in parameter space. We first estimate the volume of NROY space relative to
702 the original parameter space. This is done by Monte Carlo simulation where a large
703 number of points are uniformly drawn from parameter space and the proportion
704 within NROY space recorded. After matching to the ACC strength in section 3 we
705 had ruled out over 95% of the parameter space of HadCM3. Our current NROY
706 space is just 0.7% of the original space now that we include the North Atlantic
707 SST constraint, an estimate based on 10^6 Monte Carlo samples.

708 We can investigate features of the shape of the NROY parameter space by sam-
709 pling implausibilities. By looking at 1 and 2 dimensional representations of NROY
710 space, we can assess how different parameters combine to improve the model. Fig-
711 ure 11 shows marginal density plots for 9 of the more interesting atmosphere and
712 ocean parameters. The first panel, showing the convective cloud entrainment rate
713 coefficient, `entcoef`, indicates that low values of `entcoef` are implausible, as shown
714 by Joshi et al. (2010) for HadSM3, and that there are more NROY models at
715 the upper end of its range than in the range between 2-4 that is often determined
716 to contain the best models (Sexton et al., 2011; Rowlands et al., 2012). Isopycnal
717 diffusivity in the ocean (`ah11_si`) is also very active, with values towards the top
718 end of its range favoured. The cloud droplet to rain conversion rate, `ct`, and the
719 relative humidity threshold for cloud formation, `rhcrit`, have similar profiles as do
720 the cloud droplet to rain threshold over land, `cwland`, and the boundary layer
721 cloud fraction at saturation, `eacf`. There is also a similarity in the marginal form
722 for the ocean mixed layer parameters `lamda`, the wind mixing energy scaling factor
723 and `delta_si`, the wind mixing energy decay depth, with `lamda` in particular quite
724 active.

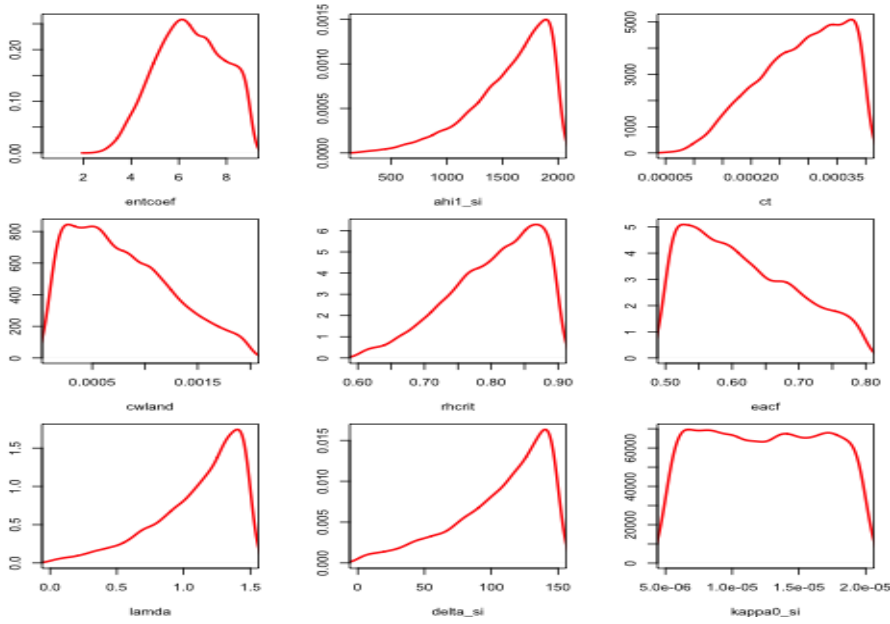


Fig. 11 Marginal NROY density plots 9 for the parameters.

725 Note that a collection of 1D marginals can be difficult to interpret. For example,
 726 the parameter setting fixed on each of the marginal modes is ruled out by SAT.
 727 More informative are the NROY density and minimum implausibility plots for 2D
 728 projections of the parameters, shown in figure 12. Each panel on the upper triangle
 729 shows the proportion of parameter settings behind each pixel that are NROY. The
 730 map is drawn by fixing the two parameters labelled for each panel at the value of
 731 a pixel, and exploring a 1000 point Latin Hypercube in the other 25 dimensions of
 732 HadCM3, plotting the proportion of samples in NROY space. Hence each upper
 733 triangle image can be viewed as a 2D projection of the density of NROY space.
 734 Grey regions are completely ruled out, meaning that, for any grey pixel, we were
 735 unable to find any NROY parameter setting in the other 25 dimensions for the
 736 given value of the other 2.

737 The standard version of HadCM3, which is ruled out by our history match
 738 to ACC strength, is plotted as the solid triangle in each panel. The version of
 739 HadCM3 we explored in section 3 is the circular point. From this picture we
 740 notice that though, perhaps, entcoef was a little low for the standard setting, the
 741 combination of low rhcrit and ct rules out the standard parameter setting anyway.

742 The NROY density plots reveal a great deal of non-linear structure to NROY
 743 space. We see that ocean parameters, such as the isopycnal diffusivity (ahi_si) must
 744 be varied jointly with cloud parameters such as ct, cwland and eacf in order to find
 745 NROY models. This result is much stronger than saying that tuning procedures
 746 that only vary one parameter at a time will not be successful for HadCM3, it
 747 says that one must tune the atmosphere and the ocean together. Parameters that
 748 appear to be reasonable at tuning a particular process or even the atmosphere

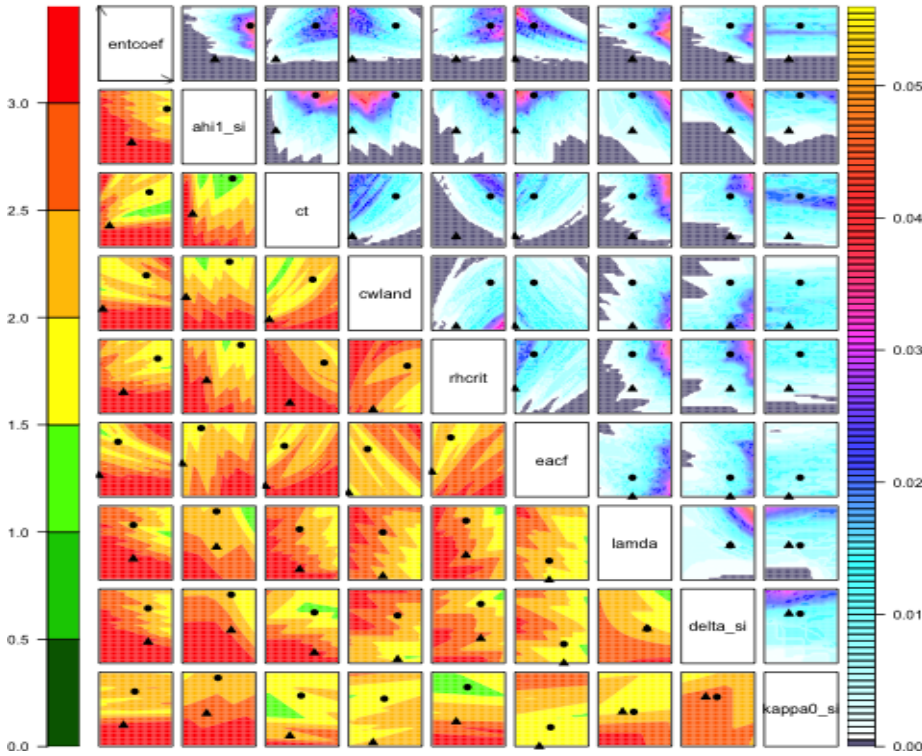


Fig. 12 NROY density plots (upper triangle) and minimum implausibility plots (lower triangle) for 2D projections of NROY space. Each panel plots either NROY density or minimum implausibility for a pair of HadCM3 parameters. NROY densities, for each pixel on any panel in the upper triangle, represent the proportion of points behind that pixel in the remaining 25 dimensions of HadCM3’s parameter space that are NROY and are indicated by the colour whose scale is indicated on the right. Minimum implausibilities, for each pixel on any panel on the lower triangle of the picture, represent the smallest implausibility found by fixing the two parameters at the plotted location and searching the other 25 dimensions of the HadCM3 parameter space. These plots are orientated the same way as those on the upper triangle, for ease of visual comparison. Standard HadCM3 is depicted on each panel as the triangular point. The parameter setting we explored in section 3 is the circular point.

749 only model, may not be close to optimal for different, and better configurations of
 750 the model ocean.

751 The lower triangle shows minimum implausibility plots. Similar to those on the
 752 upper triangle, for each pixel, representing a fixed value of a pair of parameters,
 753 a Latin Hypercube in the other 25 dimensions of HadCM3 is searched for NROY
 754 parameter values. We plot the value of the smallest implausibility found in each
 755 LHC. The plots have the same orientation as those on the upper triangle so that
 756 comparisons are easier to make. If our emulators were extremely accurate, green
 757 and yellow areas of this plot would indicate the location of potentially “good”
 758 settings of the model parameters. Though history matching restricts itself to ruling
 759 out the bad parameter settings, if that has been done and the emulators were
 760 extremely accurate (with little uncertainty in the posterior), we would look to

761 further explore the green and yellow areas of these plots as containing points that
762 are actually consistent with the data.

763 Certain panels in the lower triangle reinforce the case for varying all of the
764 parameters of the model simultaneously. Take the plot of ct against $rhcrit$ for
765 example (and remember the orientation mimics the upper triangle). Suppose we
766 keep one parameter fixed and vary the other, starting from low values of ct and
767 $rhcrit$ (as in standard HadCM3, though perhaps with a lower value of $rhcrit$ if we
768 are being strict). By doing this, according to the figure, we would never escape
769 the red zone, meaning that we would never find a model that satisfies all of our
770 metrics to within our tolerance to error. We might then come to the conclusion
771 that we had identified a structural error. But note that this figure actually implies
772 that if just one of the parameters were held fixed, and all 26 other parameters
773 varied, we would still not come across any NROY models. We would believe that
774 we had found structural errors, and yet, by varying all of the parameters, we find
775 better models (as we have) and can point to even better ones at untried parameter
776 settings.

777 The upper triangle shows that, with the exception of $entcoef$, whose values
778 cannot be low, no matter what the setting of the model parameters, values can
779 be found for each of the other parameters that would lead to a NROY model,
780 though most do not. For example, in general, the lower the isopycnal diffusivity
781 in the ocean is set, the larger the mixed layer parameter δ_{si} (which governs
782 the decay of wind mixing energy with depth) must be in order to avoid ruling
783 out that parameter choice, independent of any of the other parameters. These
784 restrictions on any given parameter range would become greater as either more
785 metrics were included or as further ensembles allowed more accurate emulators
786 to be built within parameter space (reducing the denominator in (1) and thus
787 increasing the number of models ruled out).

788 As noted previously (see figure 1), our model described in section 3 is NROY,
789 whilst the standard setting is ruled out. However, it is clear that though our
790 model may have reasonable settings for some of the parameters, there are regions
791 of parameter space with a higher density of NROY points that we would like to
792 explore, and that our model is not one of the lowest implausibility models found
793 during sampling the emulator. If we had the resource, the next step in this type of
794 analysis would be to design an ensemble within NROY space, run it, re-emulate
795 and perform another history match to further refocus the search.

796 We note that there are lots of NROY models in some of the corners of parameter
797 space. This raises questions about whether or not the pre-defined ranges of NROY
798 space were wide enough. These are valid questions which underline the importance
799 of exploring the widest physically plausible parameter ranges right from the start
800 (see discussion of this topic in Williamson et al., 2013). However, during the first
801 1 or 2 waves of a history match, the more likely explanation for seeing a lot of
802 NROY models in corners of parameter space is due to a feature of the statistical
803 modelling.

804 We design our emulators so that the uncertainty outside of the convex hull
805 of points in the ensemble (the smallest convex region containing all ensemble
806 members) increases asymptotically. This is a feature of any emulator with a mean
807 function containing non-constant functions of the parameters (Draper and Smith,
808 1998). We design our emulators in this way in order to avoid extrapolation issues
809 whereby the emulator reverts back to the prior mean outside the convex hull of

810 explored parameter settings, but with low uncertainty so that we might incorrectly
811 rule out points on the edges and in the corners. HadCM3 has 2^{27} corners and our
812 ensemble is far smaller than this, so we have many unexplored corners in parameter
813 space. This is one of the reasons we are so careful with our language. Those high
814 density regions in the corners are Not Ruled Out Yet, but it is likely that they will
815 be once our emulators can be better tuned there. If there were still high density
816 regions in corners or on edges after multiple waves, we might suspect that there
817 really were good models on the edges, and that would give us cause to consider
818 parameter values beyond the current boundary.

819 Note that it may seem natural to ask about the behaviour of models (or even
820 the model) with the lowest values of implausibility. At this stage of the analy-
821 sis, that would not be appropriate. The careful language and philosophy behind
822 history matching supports only the notion that models with large implausibilities
823 can be ruled out and says nothing about those with small implausibilities (hence
824 "not ruled out yet" as opposed to "acceptable" or even "likely"). The reason is
825 the emulator. During these early waves of analysis, much of the denominator in
826 the implausibility calculation is emulator variance, so low implausibility can mean
827 that we are simply uncertain regarding what will turn out to be a poorly match-
828 ing model. We advocate history matching as a method for tuning only if used in
829 multiple waves as discussed in section 2.1. The waves allow emulator uncertainty
830 to be reduced to the point where NROY models are close to the data (or there are
831 no NROY models). They also allow more complex joint emulators to be built more
832 easily facilitating approaches to finding the "best" models through either multi-
833 variate implausibility or otherwise (e.g. calibration as in (Edwards et al., 2011)).
834 We provide a large sample of NROY parameter values and their implausibilities
835 in the supplementary material for readers interested in further exploration of this
836 NROY space. We also include the experiment IDs and parameter values for our
837 NROY ensemble members in the supplementary material.

838 5 Discussion

839 Tuning a climate model with a high dimensional parameter space and a long run
840 time is a difficult task. Currently this task is undertaken without taking advantage
841 of the latest statistical technologies for managing uncertainty in complex models.
842 These methods allow for a targeted and comprehensive search of parameter space
843 for models satisfying numerous criteria. We have argued that many perceived
844 structural errors or "known biases" in climate models may be present due to an
845 inefficient search of the existing model parameter space during model development.

846 We have presented history matching, a technique already used to quantify
847 parametric uncertainty with climate models, as a method for climate model tun-
848 ing based on sound principles of statistical design. The method seeks to tune all
849 parameters simultaneously by using PPEs and emulators to rule out regions of
850 parameter space that lead to models that do not satisfy observational constraints
851 imposed by the model developers. We describe how the procedure should be un-
852 dertaken iteratively, with new constraints and new PPEs used to refine the search
853 for models without perceived structural biases. We also discuss how to use coarse
854 resolution versions of an expensive model so that history matching can be used to

855 assist in tuning the expensive version without the requirement for large ensembles
856 or long runs.

857 We have illustrated the power of this technique in investigating perceived struc-
858 tural biases in the HadCM3 ocean circulation. We found that the perceived struc-
859 tural bias in the ACC strength could be corrected by jointly varying both atmo-
860 sphere and ocean model parameters and showed that these changes also improved
861 other important physical properties of the ocean circulation, without compromis-
862 ing the surface air temperature profile.

863 We showed that the location of the sub-polar gyre was more realistic in the
864 model we found than in the standard HadCM3 and that the western boundary
865 current intensification in the sub-tropical gyre was greatly improved. We showed
866 that the depth profile of meridional velocities in the North Atlantic deep, unre-
867 alistic in the standard HadCM3, compares favourably with the current physical
868 understanding of these flows. We showed that the global sea surface salinity was
869 closer to observations, but that the models found in this wave of history matching
870 had a larger cold bias in the northern hemisphere SST, though the Southern Ocean
871 warm bias in the standard HadCM3 was improved. We showed that the pressure
872 and wind fields, particularly in the Southern Ocean were far more realistic in the
873 model without the ACC bias, and showed that precipitation anomalies, particu-
874 larly around the ITCZ were also improved. The AMOC and MHT are increased in
875 the improved ACC model, though the values are still consistent with observations.

876 We then illustrated the method of iterative history matching by imposing a
877 further constraint on parameter space designed to look for models without the
878 cold bias in the northern hemisphere SST. We ruled out over 99% of the model
879 parameter space as possibly containing models that satisfied our constraints and
880 showed the joint structure of the remaining space using 1 and 2 dimensional pro-
881 jections of it. We found that jointly tuning atmosphere and ocean parameters,
882 instead of tuning them one or even a few at a time was important for finding
883 these regions of parameter space. We found complex interactions between cloud,
884 ocean and convection parameters that would likely confound any one at a time
885 approach to tuning and lead us to different results than would be obtained, for
886 example, by not tuning the parameters together. For example, we found that the
887 marginal distribution for the convection parameter entcoef had much of its lower
888 range removed, yet had more density at the higher values of its range, contrary
889 to the findings of previous studies that have not simultaneously varied the ocean
890 parameters along with the atmosphere. Though this result may depend heavily on
891 our choice to tune to ACC strength, we have established through this, the need
892 for joint exploration of parameter space.

893 Though we have no further access to ensembles of HadCM3 through CPDN as
894 part of this work, further work could run an ensemble within the 1% of parameter
895 space that is NROY in order to search for even better models by history matching
896 in multiple waves. History matching is most effective with multiple waves of PPEs,
897 as the emulators improve in the region of parameter space potentially containing
898 good climates due to a higher density of model runs there. The improved emulators
899 have lower variances, which serves to increase implausibilities and rule out more
900 space.

901 Our work suggests that an overly strong ACC strength in HadCM3 is not a
902 structural error, but a calibration error. However, it may be the case that more
903 realistic ACC strengths are only possible at the expense of introducing new biases

904 in processes deemed more important than the ACC by model developers, for ex-
905 ample, it may turn out that the SST cold bias in the northern hemisphere in the
906 alternative model we studied cannot be improved without compromising the ACC
907 strength, though the results we presented in section 4.1 suggest that this would
908 not be the case. However, the best way we know of to find out for sure is to use
909 history matching with all of the important constraints included. If this is done at
910 the model development stage, structural errors in a process can be identified by an
911 attempt to history match using that process ruling out all of the parameter space.
912 If the constraints are introduced iteratively, in order of importance, the modellers
913 can determine where the structural errors are and use this information to focus
914 their research into improving the model in order to reduce or remove these errors.

915 Given the cost of developing GCMs and their importance for decision making
916 and global policy strategy, it is important that every opportunity to improve the
917 accuracy of these models is taken. History matching offers a robust and rigorous
918 statistical methodology that is easy to implement and can be used to help to
919 efficiently tune the parameters of GCMs.

920 Acknowledgements

921 This research was funded by the NERC RAPID-RAPIT project (NE/G015368/1).
922 Daniel Williamson was also funded by an EPSRC fellowship, grant number EP/K019112/1.
923 We would like to thank the CPDN team for their work on submitting our ensemble
924 to CPDN users. We'd also like to thank the Institute of Advanced Study at
925 Durham University for funding and hosting our workshop on ocean model discrep-
926 ancy which formed the motivation for these investigations. In addition, we thank
927 the oceanographers who participated in this workshop. We'd like to thank the
928 CPDN users around the world who contributed their spare computing resource as
929 part of the generation of our ensemble. NCEP and CMAP Precipitation data were
930 provided by the NOAA/OAR/ESRL PSD, Boulder, Colorado, USA, from their
931 website at <http://www.esrl.noaa.gov/psd/>

932 References

- 933 Acreman, D. M. and Jeffery, C. D. (2007), "The use of Argo for validation and tuning of mixed layer
934 models," *Ocean. Model.*, 19, 53–69.
- 935 Annan, J. D., Hargreaves, J. C., Edwards, N. R., Marsh, R. (2005), "Parameter estimation in an
936 intermediate complexity earth system model using an ensemble Kalman filter", *Ocean Modelling*,
937 8, 135–154.
- 938 Annan, J. D., Lunt, D. J., Hargreaves, J. C., Valdes, P. J. (2005), "Parameter estimation in an
939 atmospheric GCM using the ensemble Kalman filter", *Nonlinear Processes in Geophysics*, 12,
940 363–371.
- 941 Annan, J. D., Hargreaves, J. C., Ohgaito, R., Abe-Ouchi, A., Emori, S. (2005) "Efficiently con-
942 straining climate sensitivity with ensembles of paleoclimate simulations", *SOLA*, 1, 181-184,
943 doi:10.2151/sola.2005-047.
- 944 Bellprat, O., Kotlarski, S., Luthi, D., Schar, C. (2012). "Objective calibration of regional climate
945 models", *Journal of Geophysical Research*, 117, D23115.
- 946 Blaker, A. T.; Hirschi, J. J-M.; McCarthy, G.; Sinha, B.; Taws, S.; Marsh, R.; de Cuevas, B. A.; Alder-
947 son, S. G.; Coward, A. C. (2014), "Historical analogues of the recent extreme minima observed
948 in the Atlantic meridional overturning circulation at 26N". *Climate Dynamics*, 10.1007/s00382-
949 014-2274-6
- 950 Challenor, P., McNeill, D., and Gattiker, J. (2009), "Assessing the probability of rare climate
951 events," in *The handbook of applied Bayesian analysis*, eds. O'Hagan, A. and West, M., Ox-
952 ford University Press, chap. 10.

- 953 Collins, M., Brierley, C. M., MacVean, M., Booth, B. B. B. and Harris, G. R. (2007) "The Sensitivity
954 of the Rate of Transient Climate Change to Ocean Physics Perturbations," *J. Clim.*, 20, 23315–
955 2320.
- 956 Craig, P. S., Goldstein, M., Seheult, A. H., and Smith, J. A. (1996), "Bayes Linear Strategies for
957 Matching Hydrocarbon Reservoir History," in *Bayesian Statistics 5*, eds. Bernardo, J. M., Berger,
958 J. O., Dawid, A. P., and Smith, A. F. M., Oxford University Press, pp. 69–95.
- 959 Craig, P. S., Goldstein, M., Rougier, J. C., and Seheult, A. H. (2001), "Bayesian Forecasting for
960 Complex Systems using Computer Simulators," *J. Am. Stat. Assoc.*, 96, 717–729.
- 961 Cumming, J. A. and Goldstein, M. (2009), "Small sample designs for complex high-dimensional
962 models based on fast approximations," *Technometrics*, 51, 377–388.
- 963 Cunningham, S. A., Alderson, S. G., King, B. A. (2003) "Transport and variability of the Antarctic
964 Circumpolar Current in Drake Passage", *Journal of Geophysical Research*, 108, No. C5, 8084,
965 doi:10.1029/2001JC001147.
- 966 Daniel, C. (1973) "One at a time plans", *Journal of the American Statistical Association*, 68,
967 353–360.
- 968 Draper, N. R., Smith, H. (1998), "Applied Regression Analysis," 3rd Edition, John Wiley and Sons,
969 New York.
- 970 Edwards, N. R., Cameron, D., Rougier, J. C. (2011), "Precalibrating an intermediate complexity
971 climate model", *Clim. Dyn.*, 37, 1469–1482.
- 972 Fisher, R. (1926) "The arrangement of field experiments", *Journal of the Ministry of Agriculture
973 of Great Britain*, 33, 503–513.
- 974 Friedman, M., Savage, L. J. (1947) "Planning experiments seeking maxima", in *Techniques of Sta-
975 tistical Analysis*, eds Eisenhart, C., Hastay, M. W., Wallis, W. A. New York: McGraw-Hill.
- 976 Gent, P. R., Danabasoglu, G., Donner, L. J., Holland, M. M., Hunke, E. C., Jayne, S. R., Lawrence,
977 D. M., Neale, R. B., Rasch, P. J., Vertenstein, M., Worley, P. H., Yang, Z, Zhang, M. (2011),
978 "The Community Climate System Model Version 4", *Journal of Climate*, 24, 4973–4991.
- 979 Gladstone, R. M., Lee, V., Rougier, J. C., Payne, A. J., Hellmer, H., Le Brocq, A., Shepherd, A.,
980 Edwards, T. L., Gregory, J., Cornford, S. L. (2012). "Calibrated prediction of Pine Island Glacier
981 retreat during the 21st and 22nd centuries with a coupled flow line model", *Earth and Planetary
982 Science Letters*, 333-334, 191–199.
- 983 Goldstein, M and Rougier, J. C. (2009), "Reified Bayesian modelling and inference for physical
984 systems", *J. Stat. Plan. Inference*, 139, 1221–1239.
- 985 Gordon, C., Cooper, C. Senior, C. A., Banks, H., Gregory, J. M., Johns, T. C., Mitchell, J. F. B.,
986 and Wood, R. A. (2000), "The simulation of SST, sea ice extents and ocean heat transports in a
987 version of the Hadley Centre coupled model without flux adjustments," *Clim. Dyn.*, 16, 147–168.
- 988 Hargreaves, J. C., Annan, J. D., Edwards, N. R., Marsh, R. (2004), "A efficient climate forecasting
989 method using an intermediate complexity Earth System Model and the ensemble Kalman filter",
990 *Climate Dynamics*, 23, 745–760.
- 991 Haylock, R. and O'Hagan, A. (1996), "On inference for outputs of computationally expensive algo-
992 rithms with uncertainty on the inputs," in *Bayesian Statistics 5*, eds. Bernardo, J. M., Berger,
993 J. O., Dawid, A. P., and Smith, A. F. M., Oxford University Press, pp. 629–637.
- 994 Ingleby, B., and Huddleston, M., 2007: "Quality control of ocean temperature and salinity
995 profiles - historical and real-time data." *Journal of Marine Systems*, 65, 158-175,
996 doi:10.1016/j.jmarsys.2005.11.019.
- 997 Johns et al. (2006) "The New Hadley Centre Climate Model (HadGEM1): Evaluation of Coupled
998 Simulations", *Journal of Climate*, 19,1327–1353.
- 999 Johns et al. (2008) "Variability of Shallow and Deep Western Boundary Currents off the Bahamas
1000 during 2004?05: Results from the 26N RAPID?MOC Array", *Journal of Physical Oceanography*,
1001 38, 605–623.
- 1002 Joshi, M. M., Webb, M. J., Maycock, A. C., Collins, M. (2010), "Stratospheric water vapour and
1003 high climate sensitivity in a version of the HadSM3 climate model", *Atmos. Chem. Phys.*, 10,
1004 7161-7167.
- 1005 Kalnay et al. (1996), "The NCEP/NCAR 40-year reanalysis project", *Bull. Amer. Meteor. Soc.*,
1006 77, 437–470.
- 1007 Kennedy, M. C. and O'Hagan, A. (2000), "Predicting the Output from a Complex Computer Code
1008 when Fast Approximations are available," *Biometrika*, 87.
- 1009 Kennedy, M. C. and O'Hagan, A. (2001), "Bayesian Calibration of Computer Models," *J. R. Stat.
1010 Soc.. Ser. B*, 63, 425–464.
- 1011 Kraus, E. B. and Turner, J. (1967), "A one dimensional model of the seasonal thermocline II. The
1012 general theory and its consequences," *Tellus*, 19, 98106.
- 1013 Le Gratiet, L. (2014) "Bayesian analysis of hierarchical multifidelity codes", *SIAM J. Uncertainty
1014 Quantification* 1, 244-269.
- 1015 Lee, L. A., Carslaw, K. S., Pringle, K. J., Mann, G. W., Spracklen, D. V., (2011) "Emulation of
1016 a complex global aerosol model to quantify sensitivity to uncertain parameters", *Atmospheric
1017 Chemistry and Physics*, 11, pp.12253-12273. doi: 10.5194/acp-11-12253-2011.
- 1018 Loeppky, J. L., Sacks, J., Welch, W. J. (2009), "Choosing the sample size of a computer experiment:
1019 a practical guide", *Technometrics*, 51(4), 366–376.
- 1020 Martin, G. M., Milton, S. F., Senior, C. A., Brooks, M. E., Ineson, S., Reichler, T., Kim, J. (2010),
1021 "Analysis and reduction of systematic errors through a seamless approach to modelling weather
1022 and climate", *Journal of Climate*, 23, 5933–5957.

- 1023 Martin, G. M. and Coauthors (2011), “The HadGEM2 family of Met Office Unified Model Climate
1024 configurations”, *Geosci. Model Dev.*, 4, 723–757.
- 1025 Mauritsen, T., Stevens, B., Roeckner, E., Crueger, T., Esch, M., Giorgetta, M., Haak, H., Jungclaus,
1026 J., Klocke, D., Matei, D., Mikolajewicz, U., Notz, D., Pincus, R., Schmidt, H., Tomassini, L.
1027 (2012), “Tuning the climate of a global model”, *Journal of advances in modeling Earth systems*,
1028 4, M00A01, doi:10.1029/2012MS000154.
- 1029 McNeall, D. J., Challenor, P. G., Gattiker, J. R., Stone, E. J. (2013). “The potential of an observa-
1030 tional data set for calibration of a computationally expensive computer model”, *Geosci. Model
1031 Dev.* 6 1715–1728.
- 1032 Meehl, G. A., Covey, C., Delworth, T., Latif, M., McAvaney, B., Mitchell, J. F. B., Stouffer, R. J.,
1033 and Taylor, K. E. (2007), “The WCRP CMIP3 multi-model dataset: a new era in climate change
1034 research,” *Bull. Am. Meteorol. Soc.*, 88, 1383–1394.
- 1035 Megann, A. P. et al. (2010) “The Sensitivity of a Coupled Climate Model to Its Ocean Component”,
1036 *Journal of Climate*, 23, 5126–5150, doi: 10.1175/2010JCLI3394.1.
- 1037 Meijers, A. J. S., Shuckburgh, E., Bruneau, N., Sallee, J. B., Bracegirdle, T. J., Wang, Z.
1038 (2012) “Representation of the Atarctic Circumpolar Current in the CMIP5 climate models
1039 and future changes under warming scenarios”. *Journal of Geophysical Research*, 117, C12008,
1040 doi:10.1029/2012JC008412.
- 1041 Meinen, C. S. and Garzoli, S. L. (2014) “Attribution of Deep Western Boundary Current variability
1042 at 26.5 N”. *Deep Sea Research I*, 90, 81–90, <http://dx.doi.org/10.1016/j.dsr.2014.04.016>.
- 1043 Murphy, J. M., Sexton, D. M. H., Jenkins, G. J., Booth, B. B. B., Brown, C. C., Clark, R. T.,
1044 Collins, M., Harris, G. R., Kendon, E. J., Betts, R. A., Brown, S. J., Humphrey, K. A., McCarthy,
1045 M. P., McDonald, R. E., Stephens, A., Wallace, C., Warren, R., Wilby, R., Wood, R. (2009),
1046 “UK Climate Projections Science Report: Climate change projections.” *Met Office Hadley Centre*,
1047 Exeter, UK. [http://ukclimateprojections.defra.gov.uk/images/stories/projections_pdfs/
1048 UKCP09_Projections_V2.pdf](http://ukclimateprojections.defra.gov.uk/images/stories/projections_pdfs/UKCP09_Projections_V2.pdf)
- 1049 Pope, V.D. and M.L. Gallani and P.R. Rowntree and R.A. Stratton (2000), “The impact of new
1050 physical parameterizations in the Hadley Centre climate model: HadAM3.”, *Clim. Dyn.*, 16,
1051 123–146.
- 1052 Pukelsheim, F. (1994), “The three sigma rule”, *Am. Stat.*, 48, 88–91.
- 1053 Randall, D. A. and Coauthors, (2007), “Climate models and their evaluation”. *Climate Change
1054 2007: The Physical Science Basis*, S. Solomon et al. Eds., Cambridge University Press, 589–662.
- 1055 Rougier, J. C. (2007), “Probabilistic inference for future climate using an ensemble of climate model
1056 evaluations”, *Climatic Change*, 81, 247–264.
- 1057 Rougier, J. C. (2008), “Efficient emulators for multivariate deterministic functions,” *Journal of
1058 Computational and Graphical Statistics*, 17, 827 – 843.
- 1059 Rougier, J. C., Sexton, D. M. H., Murphy, J. M., and Stainforth, D. (2009), “Emulating the sensitivity
1060 of the HadSM3 climate model using ensembles from different but related experiments,” *J. Clim.*,
1061 22, 3540–3557.
- 1062 Rougier, J. C. (2013), “?Intractable and unsolved?: some thoughts on statistical data as-
1063 simulation with uncertain static parameters” *Phil. Trans. R. Soc. A*, 371, 20120297.
1064 (doi:10.1098/rsta.2012.0297) .
- 1065 Rowlands, D. J., Frame, D. J., Ackerley, D., Aina, T., Booth, B. B. B., Christensen, C., Collins, M.,
1066 Faull, N., Forest, C. E., Grandey, B. S., Gryspeerdt, E., Highwood, E. J., Ingram, W., J., Knight,
1067 S., Lopez, A., Massey, N., McNamara, F., Meinshausen, N., Piani, C., Rosier, S., M., Sanderson,
1068 B., J., Smith, L. A., Stone, D. A., Thurston, M., Yamazaki, K., Yamazaki, Y., H., Allen, M. R.
1069 (2012), “Broad range of 2050 warming from an observationally constrained large climate model
1070 ensemble”, *Nat. Geosci.*, published online, doi:10.1038/NGEO1430.
- 1071 Russell, J. L., Stouffer, R. J., Dixon, K. W. (2006) “Intercomparisons of the Southern Ocean circula-
1072 tions in IPCC coupled model control simulations”, *Journal of Climate*, 19, 4560–4575.
- 1073 Sacks, J., Welch, W. J., Mitchell, T. J., and Wynn, H. P. (1989), “Design and Analysis of Computer
1074 Experiments,” *Stat. Sci.*, 4, 409–435.
- 1075 Santner, T. J., Williams, B. J., and Notz, W. I. (2003), *The design and analysis of computer
1076 experiments*, Springer-Verlag New York.
- 1077 Schmittner, A., Urban, N. M., Shakun, J. D., Mahowald, N. M., Clark, P. U., Bartlein, P. J., Mix,
1078 A. C., Rosell-Mele, A. (2011). “Climate sensitivity estimated from temperature reconstructions
1079 of the last glacial maximum”, *Science*, 334, 1385.
- 1080 Severijns, C. A., Hazeleger, W. (2005), “Optimizing parameters in an atmospheric general circulation
1081 model”, *Journal of Climate*, 18, 3527–3535.
- 1082 Sexton, D. M. H., J. M. Murphy, and M. Collins (2011) “Multivariate probabilistic projections using
1083 imperfect climate models part 1: outline of methodology”, *Clim. Dyn.*, doi:10.1007/s00382-011-
1084 1208-9.
- 1085 Shaffrey, L. et al. (2009) “UK-HiGEM: The New UK High Resolution Global Environment
1086 Model. Model description and basic evaluation”, *Journal of Climate*, 22 (8), 1861–1896,
1087 doi:10.1175/2008JCLI2508.1.
- 1088 Solomon, S., Qin, D., Manning, M., Chen, Z., Marquis, M., Averyt, K. B., Tignor, M., and Miller,
1089 H. L. (eds.) (2007), *Contribution of Working Group I to the Fourth Assessment Report of the
1090 Intergovernmental Panel on Climate Change, 2007*, Cambridge University Press.
- 1091 Uppala et al. (2005), “The ERA-40 re-analysis”, *Q. J. R. Meteorol. Soc.*, 131, 2961–3012.

- 1092 Vernon, I., Goldstein, M., and Bower, R. G. (2010), “Galaxy Formation: a Bayesian Uncertainty
1093 Analysis,” *Bayesian Anal.* 5(4), 619–846, with Discussion.
- 1094 Watanabe, M., Suzuki, T., O’Ishi, R., Komuro, Y., Watanabe, S., Emori, S., Takemura, T., Chikira,
1095 M., Ogura, T., Sekiguchi, M., Takata, K., Yamazaki, D., Yokohata, T., Nozawa, T., Hasumi,
1096 H., Tatebe, H., Kimoto, M. (2010), “Improved climate simulation by MIROC5: Mean states,
1097 variability and climate sensitivity”, *Journal of Climate*, 23, 6312-6335.
- 1098 Williamson, D. (2010), “Policy making using computer simulators for complex physical systems;
1099 Bayesian decision support for the development of adaptive strategies,” Ph.D. thesis, Durham
1100 University.
- 1101 Williamson, D., Goldstein, M. and Blaker, A. (2012), “Fast Linked Analyses for Scenario based
1102 Hierarchies,” *J. R. Stat. Soc. Ser. C*, 61(5), 665–692.
- 1103 Williamson, D., Goldstein, M., Allison, L., Blaker, A., Challenor, P. Jackson, L., Yamazaki, K.,
1104 (2012b), “History matching for exploring and reducing climate model parameter space using
1105 observations and a large perturbed physics ensemble”, *Climate Dynamics*, To Appear.
- 1106 Williamson, D., Blaker, A. T. (2014) “Evolving Bayesian emulators for structurally chaotic time
1107 series with application to large climate models”, *SIAM/ASA J. Uncertainty Quantification*,
1108 2(1) 1-28.
- 1109 Williamson, D., Goldstein, M. (2013) “On the use of evolving computer models in Bayesian decision
1110 support”, *Journal of Statistical Planning and Inference*, In Submission.
- 1111 Williamson, D., Vernon, I. R. (2013) “Implausibility driven Evolutionary Monte Carlo for efficient
1112 generation of uniform and optimal designs for multi-wave computer experiments”, *Journal of the
1113 American Statistical Association*, In Submission.
- 1114 Xie, P., and Arkin, P. A. (1997), “Global precipitation: A 17-year monthly analysis based on gauge
1115 observations, satellite estimates, and numerical model outputs.”, *Bull. Amer. Meteor. Soc.*, 78,
1116 2539–2558.
- 1117 Yamazaki, K., Rowlands, D. J., Aina, T., Blaker, A., Bowery, A., Massey, N., Miller, J., Rye, C., Tett,
1118 S. F. B., Williamson, D., Yamazaki, Y. H., Allen, M. R. (2012), “Obtaining diverse behaviours
1119 in a climate model without the use of flux adjustments”, *Journal of Geophysical Research -
1120 Atmospheres*, Accepted.

1121 A Building emulators for history matching

1122 What follows is a brief description of the methods we used to construct emulators for the
1123 constraints described in this paper. An emulator for element i of $f(x)$ might typically be fitted
1124 as

$$f_i(x) = \sum_j \beta_{ij} g_j(x) + \epsilon_i(x) \quad (2)$$

1125 where $g(x)$ is a vector of specified functions of x , β is a matrix of coefficients, and $\epsilon(x)$ is
1126 a stochastic process with a specified covariance function. As discussed in section 2 there are
1127 many ways to build emulators and the way that is chosen will depend on the size of the PPE
1128 available, the type of constraint we wish to emulate and the relationships between the data
1129 and the parameters that we find. In this study we had access to large ensembles, and each of
1130 our constraints was a univariate quantity and so required less sophisticated modelling than a
1131 spatial field or time series might. Hence we fit the emulator mean functions, $\beta g(x)$ in equation
1132 (2) using a stepwise regression procedure described below.

1133 The functions we consider adding to $g(x)$ were linear, quadratic and cubic terms in each of
1134 the parameters with up to third order interactions between all parameters considered. Switch
1135 parameters were treated as factors (variables with a small number of distinct possible “levels”)
1136 and interactions between factors and all continuous parameters were permitted. For a list of
1137 the parameters varied in the ensemble see appendix B.

1138 Our fitting procedure begins with a “forward selection”, where we permit each allowed
1139 term to be added to $g(x)$ in its lowest available form. For example, if the linear term for x_1
1140 is not yet in $g(x)$, x_1 is available for selection but x_1^2 is not. If x_1 is already in $g(x)$ then all
1141 first order interactions with the other linear parameters in $g(x)$ are included and then x_1^2 is
1142 available for selection. So, suppose $g(x)$ is $(1, x_2)$, then the selection of x_1 implies that $g(x)$ will
1143 become $(1, x_1, x_2, x_1 * x_2)$. If x_1 is selected, at the next iteration we may select any of the other
1144 parameters but we may also include quadratic terms x_2^2 and x_1^2 . We add the interactions in this
1145 way, and do similar for third order interactions when quadratic terms have been included, so
1146 that the resulting emulator will be robust to changes of scale (see Draper and Smith, 1998,
1147 for discussion). The term that is added to $g(x)$ at each iteration is the term of those available
1148 that reduces the residual sum of squares the most after fitting by ordinary least squares.

1149 When it becomes clear that adding more terms is not improving the predictive power
 1150 of the emulator (a judgement made by the analyst based on looking at the proportion of
 1151 variability explained by the emulator and at plots of the residuals from the fit) we begin a
 1152 backwards elimination algorithm. This removes terms from $g(x)$, strictly one at a time, with
 1153 the least contribution to the sum of squares explained by the fit without compromising the
 1154 quality of the fit. Lower order terms are not permitted to be removed from $g(x)$ whilst higher
 1155 order terms remain. We stop when removing the next term chosen by the algorithm leads to a
 1156 poorer statistical model. For more details on stepwise methods such as these see Draper and
 1157 Smith (1998).

1158 We allow $\epsilon(x)$ in equation (2) to be mean zero error with variance specified by the residual
 1159 variability from the fits and no correlation between $\epsilon(x)$ and $\epsilon(x')$ for $x \neq x'$. Though this
 1160 lack of correlation might not be appropriate if we had smaller ensembles or, perhaps, if we had
 1161 completed a number of waves of history matching and were focussing on a densely sampled
 1162 subset of parameter space, it is computationally efficient and a reasonable enough approxi-
 1163 mation to the data here to be adopted for pragmatism. Including a more complex correlation
 1164 would reduce our emulator uncertainty and likely lead to more parameter space being ruled
 1165 out, though at a computational cost. Note that, with zero correlation between any points, the
 1166 emulator will not interpolate the ensemble members and will have non-zero variance at each
 1167 of them. Though the model is deterministic (in the sense that running it twice for the same
 1168 values of the model inputs returns the same answer), it also displays sensitive dependence to
 1169 initial conditions, hence the fitted variance at the design points represents the model's inter-
 1170 nal variability. This form of emulator effectively assumes that internal variability is constant
 1171 throughout parameter space, and hence can be estimated from the ensemble as part of the
 1172 fitting procedure.

1173 Following the fitting of each emulator we validate its quality using 10% of the ensemble
 1174 that was chosen randomly and reserved from the training data prior to the fit. This procedure
 1175 involves checking that the emulator accurately predicts each of the unseen ensemble members
 1176 to within the accuracy specified by emulator uncertainty. If the emulators pass this diagnostic
 1177 check, we then use them in our history matching.

Table A.1 A table indicating which terms are in $g(x)$ for our emulator of ACC in equation (2). The column and row labels refer to the parameter names as given by the matching labels in table B.1. The upper triangle labels which interaction pairs are present. The diagonal indicates the order of the highest order term in that variable. The lower triangle indicates which three way interactions are included.

	f	b	d	e	c	a	u	j	p	g	v	w	t	s	l	k
f	2	1	0	1	1	1	1	0	0	0	0	1	0	1	1	1
b		1	1	0	0	1	0	0	1	1	0	0	0	0	0	0
d			1	1	0	0	0	0	1	1	1	0	0	0	0	0
e	f			1	1	0	0	0	0	1	0	0	1	0	0	0
c					1	0	1	0	0	0	0	0	0	0	0	0
a						1	0	0	0	0	0	0	0	0	1	0
u	f						u	2	0	0	1	0	0	1	0	0
j								1	1	0	0	1	0	1	0	1
p									u	1	1	0	0	0	0	0
g										1	0	0	0	0	0	0
v											1	1	0	1	0	0
w	f											1	0	1	0	0
t													1	0	0	0
s	a													1	0	0
l															1	0
k																1

1178 We give details of our emulator for the ACC strength in HadCM3 to illustrate the complex-
 1179 ity of the mean function and the performance of the predictions. The terms selected in $g(x)$ are
 1180 displayed in table A.1. Each header corresponds to the label given to each of the parameters in

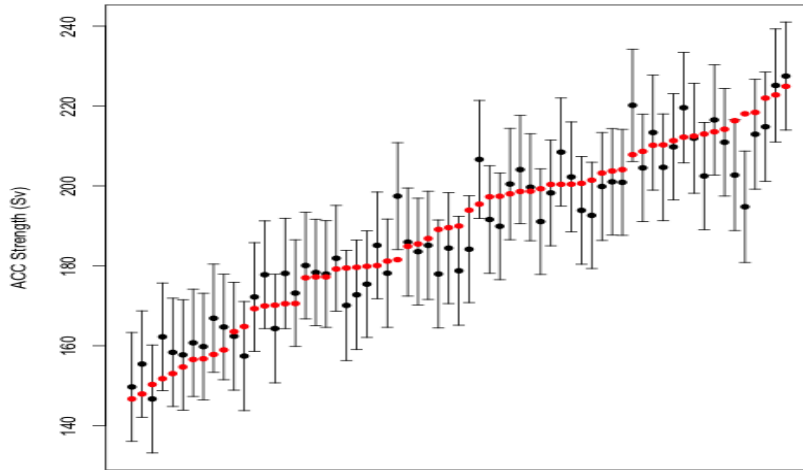


Fig. 13 Predicted ACC strength (black points) with error bars showing approximately 2 standard deviations of the emulator uncertainty for each of the withheld PPE members (red points).

1181 table B.1. Numbers on the diagonal of the table refer to the order of the parameter included in
 1182 the emulator. For example, the number 2 implies that both quadratic and linear terms in that
 1183 parameter were included in $g(x)$. Numbers on the upper triangle refer to the inclusion (1) or
 1184 not (0) of interactions between the two relevant parameters in $g(x)$. So, reading from the first
 1185 row of the table, the term (entcoef * ct) is included in $g(x)$, but the term (entcoef * rhcrit)
 1186 is not. Variables in bold on the lower triangle indicate the inclusion of the given third order
 1187 interaction. For example, the table indicates that the term (AH1.SI²*SWland) is in $g(x)$. In
 1188 addition to the terms in the table, the factor r_layers and a linear term in parameter charnock
 1189 are included, as is 1 so that an intercept is fitted.

1190 Figure 13 shows a validation plot for the ACC emulator. For 65 PPE members, chosen
 1191 randomly, that were reserved from the emulator at the fitting stage, the data are sorted by
 1192 ACC strength and plotted in red. We overlay the emulator predictions (black points) and
 1193 the uncertainty on those predictions (error bars). The uncertainty represents approximately 2
 1194 standard deviations for each prediction. We can see that the predictions are generally good
 1195 with most unseen PPE members laying within the uncertainty on the prediction. In fact, our
 1196 uncertainty specification may be too conservative, in that we have allowed for more uncertainty
 1197 in the predictions than is required. If this is the case, that would lead to less space ruled out
 1198 by history matching, not more, and it is our preference to remain conservative when ruling out
 1199 regions of parameter space.

1200 B Tables

1201 Tables B.1, B.2 and B.3 give descriptions and ranges for the parameters and the settings of
 1202 switches used in our ensemble. Some parameters have relationships with other model param-
 1203 eters that were given to us by the Met Office so that a change in one leads to a derivable
 1204 value for the other. CWland also determines CWsea, the cloud droplet to rain threshold over
 1205 sea (kg/m^3), MinSIA also determines dtice (the ocean ice diffusion coefficient) and k_gwd also
 1206 determines kay_lee_gwave (the trapped lee wave constant for surface gravity waves $\text{m}^{3/2}$).

Table B.1 Parameter descriptions and model section. CWland determines CWsea, MinSIA determines dtice and k_gwd determines kay_lee_gwave. The label column represents the labels that represent the parameters in table A.1.

Parameter	Description	Section	Label
vf1	Ice fall speed (m/s)	Cloud	a
ct	Cloud droplet to rain conversion rate (/s)	Cloud	b
CWland	Cloud droplet to rain threshold over land (kg/m ³)	Cloud	c
CWsea	Cloud droplet to rain threshold over sea (kg/m ³)	Cloud	
RHCrit	Relative humidity threshold for cloud formation	Cloud	d
eacfb1	Boundary layer cloud fraction at saturation	Cloud	e
entcoef	Convective cloud entrainment rate coefficient	Convection	f
MinSIA	Albedo at ice melting point	Sea Ice	g
dtice	Ocean ice diffusion coefficient	Sea Ice	
Icesize	Ice particle size (μm)	Radiation	h
k_gwd	Surface gravity wavelength (m)	Dynamics	i
lay_lee_gwave	Trapped lee wave constant for surface gravity waves (m ^{3/2})	Dynamics	
start_level_gwdrag	First level for gravity wave drag	Dynamics	
dyndiff	Diffusion e-folding time (hours)	Dynamics	j
dyndel	Order of diffusion operator	Dynamics	k
asym_lambda	Asymptotic neutral mixing length parameter	Boundary	l
charnock	Charnock constant	Boundary	o
cnv_rl	Free convective roughness length over sea (m)	Boundary	
flux_g0	Boundary layer flux profile parameter	Boundary	p
r_layers	No. of soil levels for evaporation	Land Surface	q
L	SO ₂ wet scavenging rate (/s)	Sulphur Cycle	
volasca	Scaling for volcanic SO ₂ emissions	Sulphur Cycle	
anthasca	Scaling for anthropogenic SO ₂ emissions	Sulphur Cycle	
so2_high_level	Model level for SO ₂ emissions	Sulphur Cycle	
vb	Background vertical viscosity (m ² /s)	Ocean	r
kb	Background vertical diffusivity (m ² /s)	Ocean	s
dkb/dz	Background vertical diffusivity gradient (m/s)	Ocean	t
AH1_SI	Isopycnal diffusivity (m ² /s)	Ocean	u
lambda	Wind mixing energy scaling factor	Ocean	v
delta_si	Wind mixing energy decay depth (m)	Ocean	w

1207 C Another NROY ACC model

1208 In the main text we present the behaviour of one of the NROY ACC models, arguing that cor-
 1209 recting the ACC strength seems to improve the ocean circulation. Though we do not reproduce
 1210 all of the figures from the main text, in order to save space, we show the BSF of another of
 1211 these models in figure 14 to indicate that the chosen model was not a “one-off”. This model has
 1212 a slightly more physical looking sub polar gyre, at the expense of a more diffuse gulf stream.
 1213 The cold bias in the North Atlantic (not shown) was also greater in this model.

1214 D Anomalies from two additional precipitation climatologies

1215 In the main text we present precipitation anomalies from the ERA40 climatology. However,
 1216 we caution the reader against interpreting improvements seen in the improved ACC run as
 1217 robust. We present two alternative precipitation climatologies, the CPC Merged Analysis of
 1218 Precipitation (CMAP) (Xie and Arkin, 1997) and NCEP/NCAR reanalysis (Kalnay et al.,
 1219 1996) in figure 15. These plots indicate that the standard model does have a tendency towards
 1220 higher than observed precipitation along the ITCZ and over the maritime continent, and the
 1221 improved ACC run we present tends to exhibit lower than observed precipitation over the
 1222 western equatorial Pacific. The improved ACC run does perform better than the standard run

Table B.2 Ranges for each of the continuous parameters varied in our ensemble. * indicates that we don't change the standard range in the exploratory sub ensemble. We don't give values for dependent parameters CWsea, dtice and kay_lee_gwave as these are calculated from CWland, MinSIA and k_gwd respectively via a one to one mapping.

Parameter	Section	Standard lower	Standard higher	New lower	New Higher
vf1	Cloud	0.5	2	0.15	2.35
ct	Cloud	5×10^{-05}	4×10^{-04}	*	5.625×10^{-04}
CWland	Cloud	1×10^{-04}	2×10^{-03}	*	*
RHCrit	Cloud	0.6	0.9	*	*
eacfb1	Cloud	0.5	0.8	*	*
entcoef	Convection	0.6	9	*	*
MinSIA	Sea Ice	0.5	0.65	*	*
Icesize	Radiation	2.5×10^{-05}	4×10^{-05}	2×10^{-05}	8×10^{-05}
k_gwd	Dynamics	$1 \times 10^{+04}$	$2 \times 10^{+04}$	*	*
dyndiff	Dynamics	6	24	*	*
asym_lambda	Boundary	0.05	0.5	0.01	0.61
charnock	Boundary	0.012	0.02	0.012	0.024
cnv_rl	Boundary	2×10^{-04}	5×10^{-03}	2×10^{-04}	6.2×10^{-03}
flux_g0	Boundary	5	20	2.5	22.5
L	Sulphur Cycle	0.33	0.33	*	*
volzca	Sulphur Cycle	1	3	0.5	3.5
anthzca	Sulphur Cycle	0.5	1.5	0.25	1.75
vb	Ocean	5×10^{-06}	8×10^{-05}	1×10^{-06}	1.1×10^{-04}
kb	Ocean	5×10^{-06}	2×10^{-05}	1×10^{-06}	3.1×10^{-05}
AH1_SI	Ocean	200	2000	100	2500
dkb/dz	Ocean	7×10^{-09}	9.8×10^{-08}	*	*

Table B.3 Switch parameters and their settings in our ensemble. * indicates that there are only 2 settings of a switch.

Parameter	Section	Setting 1	Setting 2	Setting 3
so2_high_level	Sulphur Cycle	3	5	*
start_level_gwdrag	Dynamics	3	4	5
r_layers	Land Surface	[2,1]	[3,2]	[4,3]
dyndel	Dynamics	4	6	*
lamda/delta_si	Mixed Layer	[0.3,100]	[0.5,50]	[0.7,100]

1223 everywhere outside the tropics, where the standard run has a tendency towards higher than
1224 observed precipitation.

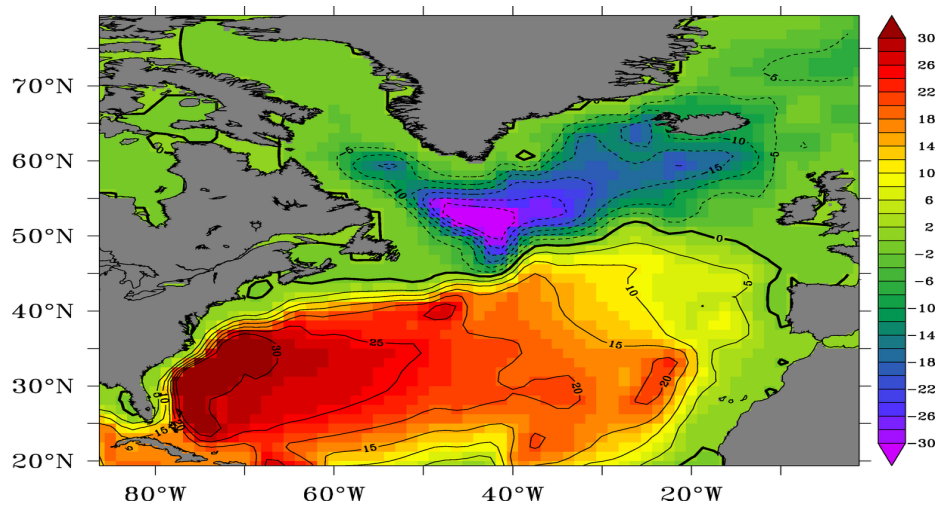


Fig. 14 The barotropic streamfunction (BSF, Sv) for a different ensemble member with realistic ACC strength (another of the blue dots from figure 1).

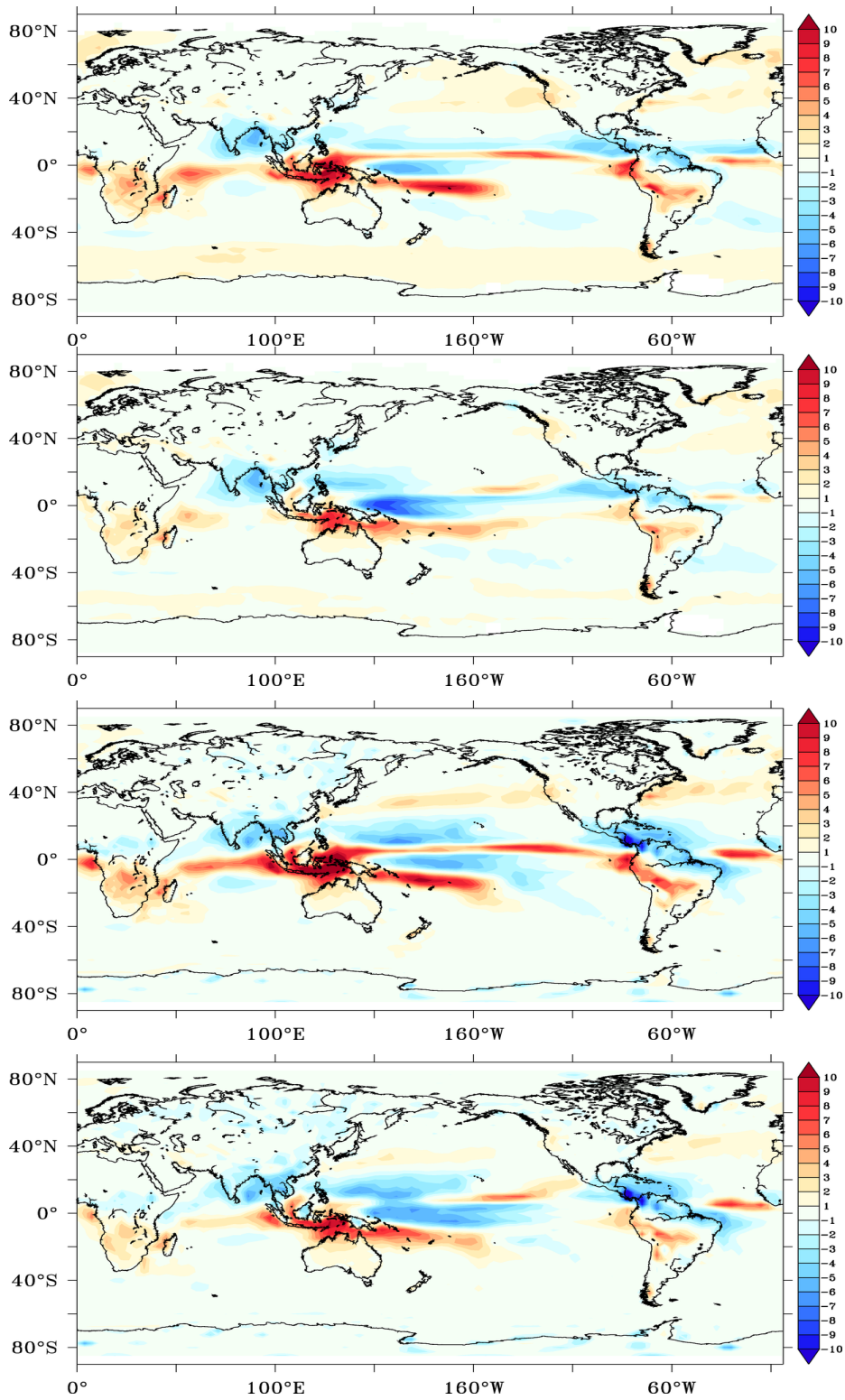


Fig. 15 The top two panels show differences from the CPC Merged Analysis of Precipitation (CMAP) climatology (standard above improved ACC model). The bottom two panels show differences from the NCEP reanalysis precipitation climatology (standard above improved ACC model).