

1 **The interpretation and use of biases in decadal climate**  
2 **predictions**

3 ED HAWKINS\*, BUWEN DONG, JON ROBSON, AND ROWAN SUTTON

*NCAS-Climate, Department of Meteorology, University of Reading, UK*

4 DOUG SMITH

*Met Office Hadley Centre, Exeter, UK*

Accepted by Journal of Climate

March 19, 2014

---

\* *Corresponding author address:* Ed Hawkins, Department of Meteorology, University of Reading, Reading. RG6 6BB. UK. E-mail: e.hawkins@reading.ac.uk

## ABSTRACT

Decadal climate predictions exhibit large biases, which are often subtracted and forgotten. However, understanding the causes of bias is essential to guide efforts to improve prediction systems, and may offer additional benefits. Here we investigate the origins of biases in decadal predictions, and whether analysis of these biases might provide useful information. We focus especially on the lead time dependent bias tendency. We initially develop a 'toy' model of a prediction system and use it to show that there are several distinct contributions to bias tendency. Contributions from sampling of internal variability and a start-time dependent forcing bias can be estimated and removed to obtain a much improved estimate of the true bias tendency, which can provide information about errors in the underlying model and/or errors in the specification of forcings. We argue that it is the true bias tendency, not the total bias tendency, that should be used to adjust decadal forecasts.

We apply the methods developed to decadal hindcasts of global mean temperature made using the HadCM3 climate model, and find that this model exhibits a small positive bias tendency in the ensemble mean. When considering different model versions we show that the true bias tendency is very highly correlated with both the Transient Climate Response (TCR) and non-greenhouse gas forcing trends, and can therefore be used to obtain observationally constrained estimates of these relevant physical quantities.

## 1. Introduction

Until recently, projections of future climate have been generated by running climate models forced by estimates of future natural and anthropogenic (e.g. from greenhouse gases and aerosols) radiative forcing. The motivation for decadal climate predictions is to improve on these standard projections by using observations to initialise predictable modes of natural variability, and by correcting errors in a model's response to past radiative forcings. Producing climate predictions that are initialised using observations of the current climate state is now a major field of scientific research (e.g., Smith et al. 2007; Keenlyside et al. 2008; Pohlmann et al. 2009; Smith et al. 2013). For example, initialised decadal climate prediction experiments are a major component of phase 5 of the Coupled Model Intercomparison Project (CMIP5; Meehl et al. 2009; Taylor et al. 2012; Meehl et al. 2013). Decadal climate predictions could potentially be of great benefit to society, for example helping to inform decisions on adaptation to a changing climate. However, there are many challenges in producing forecasts that are useful for adaptation decisions (e.g. Meehl et al. 2009; Oreskes et al. 2010).

One key challenge in producing robust predictions of future climate is to demonstrate an ability to make predictions in the past ('hindcasts'). Comparisons between hindcasts and past observations offer a wealth of information for assessing the strengths and weaknesses of

41 a prediction system, including information that can guide work to improve the system. Such  
42 an approach has proved invaluable in weather forecasting (e.g. Ferranti and Viterbo 2006).  
43 Comparisons may focus on specific case studies (e.g. Robson et al. 2012, Yeager et al. 2012),  
44 particular regions (e.g. Toniazzo and Woolnough 2013) or on the average behaviour of a  
45 system over a longer period (e.g. Smith et al. 2007, 2010; van Oldenborgh et al. 2012). A  
46 particularly important issue for decadal climate predictions is the existence of large *biases*,  
47 i.e. systematic differences between hindcasts and observations. Biases may vary with the  
48 lead time of hindcasts and are often larger than the anomalies that the system is aiming  
49 to predict. In this situation the current standard approach (e.g. Goddard et al. 2013) is  
50 to subtract the mean bias from all hindcasts before assessing other aspects of the system  
51 performance (e.g. RMSE). Such an approach is pragmatic but assumes a linear additivity  
52 between bias and forced response and ignores many important issues, such as: Why is the  
53 bias present? Does it provide any useful information? Could it be reduced?

54 The aim of this paper is to investigate the first two of these questions in particular,  
55 initially in the context of an idealised “toy” model, and secondly using results from a real  
56 decadal prediction system. We focus especially on the growth of bias with lead time, which  
57 we demonstrate offers valuable information about a prediction system and the underlying  
58 climate model. We then show further that analysis of biases for different model versions  
59 can be used to obtain useful information about the real world, in particular new constraints  
60 on the Transient Climate Response (TCR), which measures the transient sensitivity of the  
61 climate system to increases in greenhouse gases.

62 The structure of the paper is as follows. Section 2 discusses the design of decadal predic-  
63 tion experiments, and clarifies terminology. Section 3 introduces our toy model of a decadal  
64 prediction system, explains how the bias can be decomposed into distinct contributions,  
65 and examines sampling issues. The methodology we develop is then applied to predictions  
66 of global mean surface air temperature from an operational decadal prediction system in  
67 Sections 4 and 5. Conclusions and discussion of implications are in Section 6.

## 68 **2. Experimental design and terminology**

69 There are several types of decadal climate prediction experiment discussed in the litera-  
70 ture. One important issue is the specification of external radiative forcings in the hindcasts.  
71 The two main choices are:

- 72 • ‘Projection’-type - Anthropogenic forcings are assumed to be known, but ‘projected’  
73 natural forcings are used (e.g. see Smith et al. 2007). In this case any volcanic aerosol  
74 present at the forecast start time is allowed to decay, but no ‘future’ volcanic aerosol  
75 is used. In addition, the solar cycle is repeated from the previous cycle. This approach

76 attempts to mimic the realistic situation in which there is little knowledge of future  
77 natural forcing.

- 78 • 'CMIP5'-type - All forcings are assumed to be known. This is the design adopted by  
79 the CMIP5 protocol (Taylor et al. 2012).

80 In addition, hindcasts may be initialised using observations at the forecast start time  
81 ('Assim' - because assimilation is used to generate the initial states), or be initialised directly  
82 from a model state without the use of observations ('NoAssim').

83 The simplest case is arguably the 'NoAssim CMIP5' type, corresponding to traditional  
84 so-called "transient" climate model simulations. However, the ensemble sizes for these sim-  
85 ulations tend to be small (fewer than 5), which - as we will show - limits the robustness of  
86 the bias analysis. In this study we focus on the 'NoAssim projection' type of hindcasts, as  
87 performed by the UK Met Office (see Smith et al. 2007, hereafter S07). The Met Office used  
88 this approach to produce a very large ensemble of hindcasts with different versions of the  
89 same GCM (Smith et al. 2010), which proves to be a very useful resource for our analysis.  
90 However, in examining these hindcasts we must take into account the difference between the  
91 natural forcings used to force the model and those that occurred in the real world.

92 The reason that we focus on 'NoAssim'-type experiments is that understanding the biases  
93 in these experiments is a pre-requisite for understanding the biases in 'Assim'-type experi-  
94 ments. We demonstrate that the bias derived from 'NoAssim' experiments provides useful  
95 information, and we will be investigating applications to 'Assim'-type experiments in future  
96 work.

### 97 **3. Estimating bias in a toy model of a decadal predic-** 98 **tion system**

99 We first build a toy model of a decadal prediction system to examine some of the issues  
100 involved with estimating the bias of a real prediction system.

#### 101 *a. Bias of hindcasts*

102 Pseudo-observations,  $O(t)$ , are generated by assuming an externally forced linear trend  
103 in time, with added red noise,

$$O(t) = \tilde{O} + \alpha t + \epsilon(t), \quad (1)$$

104 where  $t$  is time,  $\tilde{O}$  is the 'observed' climatology,  $\alpha$  is the slope of the linear trend, and the  
105 red noise is denoted by  $\epsilon(t)$ .

106 We first assume that the ensemble mean of our pseudo-hindcasts ( $N$ ) for the same quan-  
 107 tity can be generally represented, for start time  $T$  and lead time  $\tau$ , by

$$N(T, \tau) = \tilde{N} + (T + \tau)\gamma \quad (2)$$

108 where  $\tilde{N}$  is the model climatology and  $\gamma$  is the modelled linear response to the external forc-  
 109 ing. If  $\alpha \neq \gamma$  then the climate model would produce a different trend from the observations  
 110 and therefore be biased. This could either be because the model is in error, or because there  
 111 is an error in the specification of the forcing (see later). This equation for  $N$  assumes that  
 112 we have an infinite ensemble of hindcasts, as there is no noise in the ensemble mean. This  
 113 assumption will be relaxed later. Note that these pseudo-hindcasts are only attempting to  
 114 predict the forced response, and not the internal variability component.

115 The bias ( $B$ ) of a prediction system is simply the mean error as a function of prediction  
 116 lead time,

$$B(\tau) = \frac{1}{L} \sum_{T=1}^L [N(T, \tau) - O(T + \tau)] \quad (3)$$

117 where  $L$  is the number of hindcast start dates and we assume that there is a decadal hindcast  
 118 ( $\tau = 1$  to 10 years) started every year between, and including,  $T = 1$  and  $T = L$ . Note  
 119 that in an operational system  $N$  and  $O$  would often represent anomalies from a particular  
 120 reference period. However, our analysis focusses on ‘bias tendency’ (defined below) which is  
 121 independent of the choice of reference period.

### 122 *b. Correcting the bias for observed variability*

123 The estimated bias defined in Eqn. 3 has two contributing factors, namely the true bias  
 124 (if  $\alpha \neq \gamma$  or  $\tilde{N} \neq \tilde{O}$ ) and a bias from an insufficient sampling of the internal variability in  
 125 the observations. Ideally, we would like to correct for this second variability contribution to  
 126 obtain the true bias.

127 Following Robson (2010), in the case of an infinite ensemble in a stationary climate  
 128 ( $\alpha = \gamma = 0$ ), the bias from Eqn. 3 would be,

$$B_{\text{stationary}}(\tau) = \frac{1}{L} \sum_{T=1}^L [\tilde{N} - \tilde{O} - \epsilon(t)], \quad (4)$$

$$= \tilde{N} - \tilde{O} - \frac{1}{L} \sum_{t=\tau}^{L+\tau} \epsilon(t), \quad (5)$$

$$= \tilde{N} - \tilde{O} + B_{\text{obsvar}}(\tau), \quad (6)$$

129 where  $t$  represents time and  $B_{\text{obsvar}}(\tau)$  is the mean of the observational anomalies used for  
 130 validation for a particular lead time  $\tau$ . An important point is that different observations are

131 used for different lead times. Thus  $B_{\text{obsvar}}(\tau)$  is an estimate of the bias due to the insufficient  
 132 sampling of the observed variability and will tend to zero as  $L$  increases leaving the true  
 133 bias,  $\tilde{N} - \tilde{O}$ .

134 For the more realistic case when the climate is not stationary, and there is a trend in the  
 135 observations ( $\alpha \neq 0$ ) then we can estimate:

$$B_{\text{obsvar}}(\tau) = -\frac{1}{L} \sum_{t=\tau}^{L+\tau} \text{detrended}[O(t)], \quad (7)$$

136 and this is the definition we adopt. In the toy model examples shown here we use a linear  
 137 detrending. When considering the real observations we performed sensitivity tests to explore  
 138 linear and quadratic detrending and the results were very similar (not shown), so assume a  
 139 linear detrending in all that follows.

140 A schematic demonstrating  $B_{\text{obsvar}}$  for different lead times is shown in Fig. 1 with pseudo-  
 141 observations in black, which include a linear trend and red noise, and some predictions (for  
 142 a non-infinite ensemble) shown in red in each panel. The grey regions indicate the area to  
 143 be integrated to give the value of  $B_{\text{obsvar}}$ , which varies with the lead time chosen, and need  
 144 not be zero, as shown in Fig. 1d.

### 145 *c. Bias tendency*

146 In this analysis we generally consider the ‘bias tendency’ ( $B'$ ) rather than the bias itself,  
 147 i.e. we use the bias relative to the bias for the mean of the first year,

$$B'(\tau) = B(\tau) - B(\tau = 1). \quad (8)$$

148 This choice is made because we want to consider the growth of bias with lead time, which  
 149 is natural for a prediction system. We do not use  $\tau = 0$  to avoid arbitrary assumptions  
 150 about defining climatological periods. Hence, this bias tendency has the desirable property  
 151 of being independent of the choice of climatology.

152 Similarly to the bias, the observed variability correction is also made into a tendency,

$$B'_{\text{obsvar}}(\tau) = B_{\text{obsvar}}(\tau) - B_{\text{obsvar}}(\tau = 1), \quad (9)$$

153 as shown in Fig. 1e, and an estimate of the underlying *true bias tendency* ( $B'_{\text{true}}$ ) is then,

$$B'_{\text{true}}(\tau) = B'(\tau) - B'_{\text{obsvar}}(\tau). \quad (10)$$

154 The nature of the bias growth may give valuable information about the physical processes  
 155 which cause prediction error, potentially allowing particular parameterisations to be targeted  
 156 for improvement, for example.

157 *d. Estimating the bias tendency in the toy hindcasts*

158 To test the bias tendency estimates described above, we first consider whether we can  
 159 estimate the true bias tendency of the toy model using various numbers of hindcast start  
 160 dates. Here, we generally assume that  $\alpha = 0.016\text{K/year}$  and that the red noise in Eqn. 1 ( $\epsilon$ )  
 161 has an AR(1) parameter,  $\beta = 0.5$ , and total variance,  $\sigma_\epsilon^2 = 0.01\text{K}$ . These values are chosen to  
 162 roughly simulate observed annual global mean surface air temperature (SAT) observations  
 163 since 1850 (Brohan et al. 2006), although the conclusions are insensitive to the exact choices.  
 164 We pick  $\gamma = 0.020\text{K/year}$ , i.e. the toy hindcasts are positively biased by 25%, and retain the  
 165 infinite ensemble assumption for now.

166 An example of such a hindcast system is shown in Fig. 2a for decadal hindcasts started  
 167 every year for  $L = 20$  years, where the black line represents the observations, the solid blue  
 168 line is the true forced trend ( $\alpha$ ), the dashed blue line is a linear fit to the observations used  
 169 in the estimation of  $B'_{\text{obsvar}}$ , and the red lines represent the pseudo-hindcasts ( $N$ ) which are  
 170 identical because of the infinite ensemble assumption.

171 In Fig. 2b, we show estimates of the bias tendency for the situation in Fig. 2a. The solid  
 172 blue line uses the definition of uncorrected bias tendency (Eqn. 8), and the dashed blue line  
 173 corrects for the observed variability using Eqn. 10. Note that the dashed blue line does not  
 174 match the true bias (grey shading) because the estimated trend from the observations is  
 175 not correct, i.e. the estimate of  $B'_{\text{obsvar}}$  is not exact. If the true forced trend is used in the  
 176 estimation of  $B'_{\text{obsvar}}$  then the true bias tendency is recovered (black line).

177 We next simulate 1000 realisations of the pseudo-observations and hindcast sets. Bias  
 178 tendency estimates for 10 examples of these realisations are shown in Fig. 2c. With these  
 179 20 start dates there is a wide range of estimated bias tendencies. For different numbers of  
 180 hindcast start dates ( $L$ ), Fig. 3 demonstrates that correcting the bias tendency using  $B'_{\text{obsvar}}$   
 181 (dashed line) reduces the error in the estimates of bias tendency at a lead time of 10 years  
 182 compared to using the uncorrected bias tendency (solid line). Both estimators of the bias  
 183 tendency are themselves unbiased, i.e. the mean over all realisations equals the true bias  
 184 tendency (not shown). The spread in bias tendency estimates decreases with the number  
 185 of start dates as more observations allow more accurate estimates. The observed variability  
 186 correction also becomes smaller with more start dates. When analysing the operational  
 187 NoAssim hindcasts in Section 4 we generally use 40 start dates, so the spread is around half  
 188 as large as suggested in Fig. 2c.

189 For the particular set of toy model parameters chosen here, we see that the expected error  
 190 in the bias tendency estimate becomes smaller than the bias itself (grey line in Fig. 3), i.e.  
 191 the sign of the true bias tendency could be detected, for around  $L = 15 - 20$  hindcast start  
 192 dates. For fewer hindcasts, the uncertainty in the bias estimates does not allow a detection,  
 193 with the implication for ensemble design that more start dates are required. If the bias is

194 uncorrected then more start dates are required to detect the bias.

195 *e. Forcing bias and consistent verification times*

196 So far we have assumed that the radiative forcing that is causing a warming or cooling  
 197 trend has been correctly specified and so any bias tendency is due to errors in the model  
 198 response to this forcing. However, there are two types of forcing bias which could make  
 199 this assumption invalid, namely start-time independent and start-time dependent bias. The  
 200 'CMIP5' design discussed in Section 2 results in start-time independent forcing biases because  
 201 all hindcasts see the same forcing at the same date. However, for the 'Projection' design  
 202 this is not the case: hindcasts started from different dates may see different forcings. For  
 203 example, a hindcast started in 1989 would not include any volcanic aerosol from the Mount  
 204 Pinatubo eruption in 1991, whereas a hindcast started in 1992 would. Thus there is a start-  
 205 time dependent forcing bias. S07 noted that this type of forcing bias makes a significant  
 206 contribution to the bias of a set of hindcasts. They attempted to remove it, somewhat  
 207 arbitrarily, by excluding years just after volcanic eruptions from the estimation of the bias.  
 208 Fortunately, a further correction is available to account for this start-time dependent bias.

209 In deriving,  $B$  from Eqn. 3 we chose to use all possible combinations of start dates and  
 210 verification times. However, an alternative is to use a 'consistent' set of verification times,  
 211 which only includes years where all lead times,  $\tau$ , can be *simultaneously* assessed, i.e. the  
 212 *same* observation can be used to assess the bias at all lead times. In the schematic of Fig. 1  
 213 these times are shown by the range of the blue bars, i.e. years 11-21 in this example, as year  
 214 11 is the earliest time that a 10 year lead time forecast can be verified (along with forecasts  
 215 for lead times of 1-9 years), and year 21 is the last time that a 1 year lead time can be  
 216 verified (along with forecasts for lead times of 2-10 years).

217 Using these consistent verification times, assuming there is no start time dependent  
 218 forcing bias and an infinite ensemble, and generalizing from Eqn. 3, the bias becomes,

$$B_{\text{consis}}(\tau) = \frac{1}{L - \tau_{\text{max}} + 1} \sum_{t=1+\tau_{\text{max}}}^{L+1} [N(T, \tau) - O(t)], \quad (11)$$

$$= \frac{1}{L - \tau_{\text{max}} + 1} \sum_{t=1+\tau_{\text{max}}}^{L+1} [N(t) - O(t)], \quad (12)$$

$$= A, \quad (13)$$

219 where  $\tau_{\text{max}}$  is the largest lead time to be considered. Crucially, for this particular choice  
 220 of verification times, all the terms on the right hand side of Eqn. 12 are *independent* of  
 221 lead time, because  $N(t)$  is the same for all lead times and  $B'_{\text{obsvar}}$  is zero for this choice of  
 222 verification times (Fig. 1). In this instance,  $B_{\text{consis}}(\tau)$  is a constant ( $A$ ) with lead time, and

223 therefore, the bias tendency using consistent verification times is,

$$B'_{\text{consis}}(\tau) = B_{\text{consis}}(\tau) - B_{\text{consis}}(1), \quad (14)$$

$$= 0. \quad (15)$$

224 Hence, *in the absence of a start-time dependent forcing bias*,  $B'_{\text{consis}}$  is exactly zero (assuming  
225 an infinite ensemble).

226 To test the impact of a start-time dependent forcing bias in our toy model, we generalise  
227 Eqn. 1 by adding a volcanic eruption into the pseudo-observations, within the consistent  
228 validation time period, of the form,

$$V(\xi) = 0.2 \exp(-\xi) \quad (16)$$

229 where  $V$  is the temperature response to a volcanic eruption, which reduces over time ( $\xi$ ,  
230 measured in years) with an exponential decay timescale of 1 year, from a peak impact of  
231 0.2K. We also assume that the hindcasts also include this impact, but only after the eruption  
232 has occurred.

233 Repeating our toy hindcasts (Fig. 4), still assuming an infinite ensemble, demonstrates  
234 that the measured bias tendency (blue) is over-estimated when compared to the true bias  
235 tendency (dark grey), because the bias tendency due to the volcanic eruption is non-zero  
236 (light grey).  $B'_{\text{consis}}$  is shown by the red line in Fig. 4b, which matches the forcing bias  
237 tendency (light grey) as expected.

238 Note especially that to estimate  $B'_{\text{consis}}$  from the data there is no need to assume any  
239 functional form for the forcing bias. Therefore, we can correct for the start-time dependent  
240 forcing bias by estimating the bias tendency using all verification times, and subtracting  
241 off the bias tendency estimated using consistent verification times ( $B'_{\text{consis}}$ ). Generalising  
242 Eqn. 10,

$$B'_{\text{true}}(\tau) = B'(\tau) - B'_{\text{obsvar}}(\tau) - B'_{\text{consis}}(\tau). \quad (17)$$

243 The green lines in Fig. 4b are an example of such an estimate using the bias tendency  
244 corrected only by the consistent verification times (solid) and using Eqn. 17 (dashed). Below  
245 we will demonstrate that it is necessary to remove the forcing bias in this way to obtain a  
246 robust estimate of the true bias tendency, which is the key quantity of interest.

247 We note here that there are still two contributions to the true bias tendency. The first  
248 is errors in the underlying climate model; for example, if the sensitivity of the model to  
249 greenhouse gas forcing is higher or lower than that of the real world, the hindcasts will  
250 warm too rapidly or too slowly, giving a positive or negative bias tendency. The second  
251 is (start-time independent) errors in the forcing applied to the model; for example, if the  
252 negative radiative forcing due to anthropogenic aerosols is lower or higher in the model than  
253 in the real world, this will also give a positive or negative bias tendency. Correcting the

254 bias tendency using the period of consistent verification times does not deal with the issue  
 255 of forcing errors that may occur outside of the period of consistent verification times, and  
 256 this is discussed further when considering the real observations.

257 Finally, it should be noted that estimating the bias tendency using all verification times  
 258 and subtracting off the bias tendency using consistent verification times is not the same as  
 259 estimating the bias tendency using ‘non-consistent’ verification times (not shown).

260 *f. How many ensemble members are needed?*

261 As discussed above, we have so far assumed that the toy hindcasts have infinite ensemble  
 262 members. We now relax this assumption to understand how many ensemble members would  
 263 be required to ensure a robust bias tendency estimate.

264 For a finite ensemble, our toy model for the predictions is generalized from Eqn. 2 to,

$$N(T, \tau) = (T + \tau)\gamma + \zeta(T, \tau) \quad (18)$$

265 where  $\zeta$  is red noise with the same AR1 parameter as the pseudo-observations ( $\beta = 0.5$ )  
 266 and a noise component which depends on  $M$ , the number of ensemble members, i.e.  $\sigma(\zeta) =$   
 267  $\sigma_\epsilon/\sqrt{M}$ . Note that this definition is equivalent to taking the mean of  $M$  different ensemble  
 268 members, each with variance  $\sigma_\epsilon^2$ .

269 Fig. 5 explores the spread in estimates of the true bias tendency using various values for  
 270  $M$ , making (or not) the different corrections discussed above. This spread is derived from  
 271 100,000 different realisations of the toy model. The colours represent using 20 start dates  
 272 (grey) and 40 start dates (blue). Firstly, the most reliable and accurate estimate of the true  
 273 bias is when all the corrections described above are applied (Fig. 5a). For the other cases,  
 274 the bias estimate itself becomes more biased, or more uncertain (Fig. 5b,c,d).

275 In addition, as the number of ensemble members is increased the uncertainty in the bias  
 276 estimates initially decrease, but then stabilise. For  $M \gtrsim 8$ , the expected error in the bias  
 277 remains roughly constant. This analysis suggests that as long as  $M \gtrsim 8$ , then the ensemble is  
 278 effectively infinite for global mean temperature. In addition, to detect the sign of a true bias  
 279 tendency *it is far better to increase the number of start years, than to increase the number*  
 280 *of ensemble members.* This is also found to be the case when the variance of the noise is  
 281 doubled to represent a regional mean, rather than a global mean (not shown).

282 We note that the mean of the toy model realisations in the fully corrected case does not  
 283 quite match the expected value (black). This is probably due to an interaction between the  
 284  $B_{\text{consis}}$  and  $B_{\text{obsvar}}$  correction terms as  $B_{\text{consis}}$  will also have a variability component, but this  
 285 estimate is still the least biased.

## 4. Estimating the true bias in an operational decadal prediction system

S07 describe the performance of a set of hindcasts made using the HadCM3 global climate model (Gordon et al. 2000). Here we analyse a later set of ensembles, termed NoAssimPPE, which utilises the same HadCM3 GCM, but with 9 different ‘perturbed physics’ versions (Smith et al. 2010). These different PPE versions were chosen to sample a wide range of climate sensitivities and ENSO amplitudes (e.g. Murphy et al. 2004; Smith et al. 2010; Collins et al. 2011).

The hindcasts were initialised from model states consistent with the applied radiative forcings using start dates once per year from 1961-2001, with one 10 year prediction per model version. As in the original S07 hindcasts, the NoAssimPPE hindcasts used the ‘Projection’ approach to specifying external forcings (Section 2).

### *a. Start-time dependent forcing bias*

First, we demonstrate the presence of a start-time dependent forcing bias in the NoAssimPPE hindcasts (41 start dates, 9 ensemble members, 1961-2001). Because the hindcasts use only information available at the start of the forecast, ‘future’ volcanic eruptions were not considered. This produces hindcasts that are biased warm when compared to observations. Also, the previous solar cycle is repeated, which is another potential source of bias.

Fig. 6 shows estimates of the natural forcings (volcanic and solar) used in the transient 20th century integrations (left) and in the prediction system (middle). The estimates for the prediction systems assume an exponential decay rate of the volcanic aerosol present at the forecast start time of 1 year and an 11 year solar cycle length. The resulting forcing bias is shown in the right column.

When integrated over all start dates an estimate of the start-time dependent forcing bias is produced (bottom right). The magnitude of the bias is dominated by the volcanic component and peaks at around  $0.45\text{Wm}^{-2}$  at a lead time of 3 years, subsequently dropping to around  $0.30\text{Wm}^{-2}$  at a lead time of 10 years.

### *b. Bias tendency estimates in NoAssimPPE*

We now explore the expected error in the bias estimates using the results from analysis of the toy model. Fig. 7 shows the expected growth with lead time of the error in the estimated bias for NoAssimPPE (grey) where the solid (dashed) grey line indicates the expected error using 1 (9) ensemble members. The black line shows the corresponding error for the original NoAssim (S07) hindcasts (effectively 20 start dates and 16 ensemble members). The greater number of ensemble members in the original NoAssim results in a smaller expected error at

320 short lead times (1-3 years), compared with the single member PPE system. However, the  
321 larger number of start dates in NoAssimPPE suggests a far smaller error at long lead times  
322 (5-10 years), even using a single ensemble member. The uncertainty estimates for 5-year  
323 means (horizontal grey bars) are used below in Section 5.

324 We next apply the bias estimate methodology developed using the toy model to annual  
325 means of global mean surface air temperature from the NoAssimPPE hindcasts (Fig. 8).  
326 We compare the hindcasts to four observational datasets (HadCRUT4 - Morice et al. 2012,  
327 GISTEMP - Hansen et al. 2010, NCEP - Kalnay et al. 1996, ERA-40 - Uppala et al. 2005), but  
328 all give consistent results. Note that the observations used are for 1961-2010, except ERA-40  
329 which uses 1961-2001. Unless otherwise stated we use HadCRUT4 in all that follows. For  
330 the NoAssimPPE system, the raw bias tendency estimate (Fig. 8a) suggests that HadCM3  
331 has a warm bias, which is apparently primarily due to a start-time dependent forcing bias  
332 (Fig. 8b) rather than an insufficient sampling of the observational variability (Fig. 8c). The  
333 best estimate for the true bias tendency (Fig. 8d) shows a very slight warm bias of around  
334 0.02K/decade, which is marginally statistically significant. The interpretation of this true  
335 bias tendency is discussed in Section 5.

336 In addition, we note that the bias is positive over both land and sea (Fig. 8e,f). Both  
337 the spatial pattern and physical processes responsible for the bias growth will be explored  
338 in future work.

339 The global mean SAT bias tendency associated with the time dependent forcing error  
340 makes the largest contribution to the SAT total bias tendency (Fig. 8). Smith et al. (2007)  
341 also recognised the importance of accounting for the bias caused by volcanic eruptions. They  
342 estimated that the raw bias for NoAssim was around 0.14K/decade (consistent with Fig. 8),  
343 but they removed the forcing bias by excluding some years following volcanic eruptions. We  
344 believe that our result is more robust as we are accounting for the forcing bias more explicitly  
345 and objectively.

346 The lead time evolution of the ensemble mean global averaged shortwave radiation (SW)  
347 bias tendency over the ocean at the top of atmosphere (TOA) (i.e. the forcing error) using  
348 the consistent verification times is illustrated in Fig. 9a, and shows a rapid increase in  
349 downward solar radiation in the first 3-4 years to about  $0.30 - 0.35\text{Wm}^{-2}$  and it maintains  
350 this magnitude afterwards. This estimated forcing error and its lead time evolution are  
351 consistent with the implied surface heat flux bias tendency from vertically integrated ocean  
352 heat content (OHC) bias tendency (the implied flux bias tendency is not sensitive to the  
353 depth chosen for the integration since OHC bias tendency is mostly confined in the top 500  
354 metres) as shown in Fig. 9b and it is also consistent with the directly estimated forcing  
355 error associated with volcanic eruptions (Fig. 9c, smoothed from Fig. 6). A caveat with  
356 using the 1961-2001 start dates for validation is that the Agung volcano in 1963 is before the

357 consistent verification times. We have performed a sensitivity test by excluding the hindcasts  
358 from 1961, 1962, 1963, and 1964, but this does not significantly affect the results.

359 The relative importance of each component of the bias is illustrated in Fig. 10, which  
360 confirms that the lead time dependent forcing bias dominates. For NoAssimPPE the sam-  
361 pling correction (orange) is very small for global mean temperature because the number of  
362 hindcast starts dates is large. Note, however, that this contribution is expected to be larger  
363 for other variables and smaller regions. These results illustrate clearly the importance of  
364 decomposing the bias into its different components before interpreting its meaning. Fur-  
365 thermore, if a bias correction were to be applied to a *forecast* (rather than a hindcast), we  
366 suggest it is the underlying true bias tendency that should be used, rather than the raw  
367 bias tendency derived from the hindcasts, in contrast to some current practices (e.g. Smith  
368 et al. 2013). We plan to explore the issues surrounding the application of bias corrections  
369 to forecasts in future work.

## 370 5. Interpretation of the true bias tendency

### 371 a. Role of ocean heat uptake in bias tendency

372 The true bias tendency could arise either from start-time independent errors in the forc-  
373 ings applied to the model (e.g. errors in the specification of anthropogenic aerosols) or from  
374 errors in the transient sensitivity of the model to such forcings (or both). Errors in the tran-  
375 sient sensitivity could themselves arise from errors in the representation of atmospheric or  
376 surface feedbacks and/or from errors in the representation of ocean heat uptake (e.g. Raper  
377 et al. 2002, Gregory and Forster 2008, Boé et al. 2009). This last factor can be examined  
378 by considering the bias tendency for global mean ocean heat content (OHC; Figure 11).  
379 As for surface air temperature the total bias is dominated by the forcing bias. The true  
380 bias tendency for the surface or top 100m is again positive, and is near zero below a few  
381 hundred metres. If insufficient ocean heat uptake were the cause of the warming bias at the  
382 surface we would expect to see a cooling bias subsurface. The fact that we don't see such a  
383 feature suggests that ocean heat uptake is not the reason for the warming bias in surface air  
384 temperature.

385 Further insights into the true bias tendency may be obtained by considering the biases  
386 associated with individual model versions (as distinct from the ensemble mean considered  
387 previously). Figure 12 shows that, within the PPE ensemble, there is a high positive cor-  
388 relation between the true bias tendency for OHC and that for SAT. This correlation again  
389 implies that variations in ocean heat uptake are not the primary cause of variations in SAT  
390 bias in NoAssim PPE.

391 *b. Relating climate sensitivity, forcing trends and bias tendency*

392 Next we consider the possible causes of the different true bias tendencies in the various  
393 PPE versions.

394 The first possible explanation is that the true bias tendency is directly related to the  
395 climate sensitivity of the model version (Figure 13a). Values for the Transient Climate  
396 Response (TCR) were obtained for each model version through separate specific experiments  
397 carried out at the UK Met Office. The HadCM3 NoAssim PPE model versions have a TCR  
398 range of 1.6-2.7K with a mean of 2.1K, which may be compared with the likely range of 1.0-  
399 2.5K from IPCC AR5 (Stocker et al. 2013). Figure 13a shows a linear relationship between  
400 the true bias tendency for global mean SAT and TCR, in which the most sensitive models give  
401 the largest warming bias tendency, with a correlation coefficient of 0.89. This high correlation  
402 suggests that the true bias tendency may be providing very useful information about the  
403 sensitivity of the underlying model. The correlation between TCR and the uncorrected bias  
404 tendency is 0.75, so the corrections have also improved this relationship. In addition, since  
405 a perfect model should yield a true bias tendency of zero, we can use this relationship to  
406 estimate a likely range for TCR.

407 A Monte-Carlo approach is used to fit regression lines to the data by perturbing the  
408 true bias tendency of each model version, taking into account the bias tendency uncertainty  
409 (0.016K, calculated from the toy model). The distribution of the intercepts of these lines with  
410 the  $y = 0$  line (corresponding to zero true bias tendency) then provides an observationally  
411 constrained range for TCR. We find that the 5 – 95% range for TCR constrained in this way  
412 is 1.4-1.8K with a median of 1.6K using HadCRUT4 (Figure 13c). This range is considerably  
413 narrower than the corresponding likely range from IPCC AR5 of 1.0-2.5K, and observation-  
414 based ranges of 1.3-2.3K (Gregory and Forster 2008) and 0.9-2.0K (Otto et al. 2013). With  
415 doubled estimates for the uncertainty in the true bias tendency the range from this study  
416 becomes 0.9-1.9K. The standard version of HadCM3 has a TCR of 2.0K (Randall et al.  
417 2007).

418 The constrained ranges of TCR for different observational data sets, are summarized in  
419 Table 1. Results indicate that the median and the ranges of the constrained TCR are only  
420 slightly sensitive to the data that is used to validate the hindcasts, with the other datasets  
421 producing values of TCR about 0.15K higher. The reduced spread of TCR is a robust  
422 feature and so the underlying SAT true bias tendency from the decadal climate hindcasts  
423 could be used to constrain the model TCR, complementing other approaches proposed in the  
424 literature (e.g. Allen et al. 2000, Stott and Forest 2007, Gregory and Forster 2008, Knutti  
425 and Tomassini 2008, Murphy 2010, Tett et al. 2013). It is also interesting to note that  
426 having a range of models with widely different TCR has proved very useful in this analysis,  
427 especially to constrain the upper end of our TCR ranges.

428 However, there is another possible explanation for the true bias tendency differences.  
429 When considering the role of TCR we have assumed that the forcing trends in each PPE  
430 version are the same. However, Harris et al. (2013) recently demonstrated that the different  
431 PPE versions of HadCM3 have different non-greenhouse gas (GHG) forcing, likely due to  
432 the different interaction of aerosols with low clouds. The relationship is such that versions  
433 of HadCM3 with a low TCR, and negative bias tendency, also have a cooling trend from  
434 non-GHG forcing from 1961-2010, and this could potentially contribute to the relationship  
435 between TCR and true bias tendency.

436 Figure 13b relates the true bias tendency to the non-GHG forcing trends for the different  
437 PPE model versions. The forcing data are taken from Harris et al. (2013), and linear trends  
438 have been fitted from 1961-2010, excluding years with, and shortly after, volcanic eruptions.  
439 This provides an estimate of the non-GHG forcing trends and the observed relationship can  
440 be used to produce an improved constraint on the non-GHG forcing trend, which is found  
441 to be negative, unlike in the majority of the model versions.

442 There are therefore two possible causes for the relationship between perturbed parameter  
443 versions of HadCM3 and the true bias tendency, i.e. it is clear that the parameter perturba-  
444 tions affect both the TCR and the non-GHG forcing trends and that both factors influence  
445 the true bias tendency. Trying to separate the two effects is beyond the scope of this paper,  
446 but further work will use the spatial patterns, and other climate variables, to further under-  
447 stand the causes of the bias tendencies. However, we note that if both factors are playing a  
448 role then the constrained ranges for TCR and non-GHG forcing would broaden.

449 An additional related caveat is that if there is a *systematic* error (i.e. common to all  
450 model versions) in the trends in the radiative forcing applied to the model then this would  
451 also affect the true bias tendency. For example, if the forcing trends were systematically  
452 too large then the true bias tendency would also be too large, and vice versa. The result  
453 of any such bias would be to displace all the data in Figures 13a,b vertically along the true  
454 bias tendency axis. Such a displacement would shift the constrained ranges but would not  
455 broaden the distributions. This caveat should be kept in mind when interpreting our results.

456 One possible approach to addressing these various caveats would be a multi-model study  
457 where the forcings are likely to be different for each model, and this is planned further work.

## 458 6. Conclusions and discussion

459 We have explored the estimation of bias in a toy model of a decadal prediction system,  
460 and applied the techniques developed to analyse the bias of operational predictions of global  
461 mean temperature. We have focused on hindcasts initialised from model states, rather than  
462 from observations, and examined the bias tendency in particular. The main findings can be

463 summarised as follows:

- 464 • The total bias tendency can be separated into several components, namely: a contri-  
465 bution from sampling uncertainty due to internal variability, a start-time dependent  
466 forcing bias tendency, and the true bias tendency.
- 467 • We have shown how the contributions from sampling uncertainty and start-time de-  
468 pendent forcing bias can be estimated, and removed, to give a better (lower variance  
469 and less biased) estimate of the true bias tendency. We argue that it is the true bias  
470 tendency, not the total bias tendency, that should be used to adjust decadal forecasts.
- 471 • The true bias tendency is attributable to: 1) errors in the sensitivity of the underlying  
472 model to forcing, and/or 2) start-time independent errors in the specification of forcing  
473 (e.g. errors in the specification of anthropogenic aerosols).
- 474 • To improve estimates of bias tendencies, more hindcast start dates are more beneficial  
475 than more ensemble members.
- 476 • The UK Met Office NoAssim PPE prediction system exhibits, in the ensemble mean, a  
477 small positive true bias tendency in hindcasts of global mean surface air temperature,  
478 and this is marginally statistically significant. We have demonstrated that this bias is  
479 not attributable to insufficient ocean heat uptake.
- 480 • The different true bias tendencies in global mean surface air temperature in the various  
481 PPE versions can be used to constrain relevant physical properties of the models, such  
482 as the TCR and non-GHG forcing trends.

483 There are a number of caveats to the findings above. In the toy model, we have assumed  
484 linear trends. However, we do not believe that this compromises the decomposition of the  
485 bias tendency into its different terms. Secondly, we assumed that the toy model has the  
486 same variability properties as the toy observations. This is unlikely to hold perfectly in  
487 an operational setting as there is a broad spread in simulated variability amongst different  
488 models (Hawkins and Sutton 2012) and even amongst the different PPE versions of HadCM3  
489 (Ho et al. 2013), but this would only change the number of start dates and ensemble members  
490 required to reliably estimate the bias. Most importantly, we have assumed the radiative  
491 forcings imposed in the decadal hindcasts are correct, as discussed in Section 5.

492 In the decadal hindcast experiments for CMIP5, the standard start dates are every 5  
493 years (Meehl et al. 2009; Taylor et al. 2012). In this situation there is no way of estimating  
494 the consistent bias on annual timescales. Therefore, any lead time dependent errors in  
495 the forcing cannot be removed. However, in the ‘Tier 1’ CMIP5 predictions, the complete

496 volcanic and solar forcings are assumed known, so there should be little start-time dependent  
497 forcing bias. In other suggested experiments this is not the case. We suggest that the design  
498 of future decadal prediction experiments should consider start dates every year to allow for  
499 any start-time dependent forcing bias to be removed.

500 We believe that the analysis of bias tendencies has considerable potential to provide  
501 further insights into climate models and the real climate system. We note that Masson  
502 and Knutti (2013) suggest that perturbed-physics ensembles and multi-model ensembles can  
503 behave differently and show opposite emergent constraints so it would be valuable to repeat  
504 this analysis using a wider range of operational prediction systems.

505 Beyond the global means considered in this paper there is a great deal of information  
506 in the spatial patterns of bias growth for a range of variables, and we have begun work  
507 to analyse these patterns. Lastly, there is an obvious need to examine how the growth of  
508 biases in a system initialised from model states is related to the growth of biases in a system  
509 initialised from observational states. This work involves many challenges but is essential for  
510 the development of decadal predictions.

511 *Acknowledgments.*

512 We thank Glen Harris for providing the forcing data from his important study, and for  
513 valuable discussions. We also thank two anonymous reviewers for their helpful comments  
514 which improved the manuscript. The research leading to this paper has received support  
515 from NCAS-Climate (EH, BD and RS), from the European Community's 7th framework  
516 programme (FP7) under grant agreement No. GA212643 (THOR) (EH, DS) and from the  
517 UK NERC funded EQUIP (EH) and VALOR (JR) projects. DS was also supported by the  
518 joint DECC/Defra Met Office Hadley Centre Climate Programme (GA01101), and the EU  
519 FP7 COMBINE project.

520

521

## REFERENCES

522 Allen, M. R., P. A. Stott, J. F. B. Mitchell, R. Schnur, and T. L. Delworth, 2000: Quantifying  
523 the uncertainty in forecasts of anthropogenic climate change. *Nature*, **407**, 617–620, doi:  
524 10.1038/35036559.

525 Boé, J., A. Hall, and X. Qu, 2009: Deep ocean heat uptake as a major source of spread  
526 in transient climate change simulations. *Geophysical Research Letters*, **36**, L22 701, doi:  
527 10.1029/2009GL040845.

- 528 Brohan, P., J. J. Kennedy, I. Harris, S. F. B. Tett, and P. Jones, 2006: Uncertainty estimates  
529 in regional and global observed temperature changes: a new dataset from 1850. *J. Geophys.*  
530 *Res.*, **111**, D12106, doi:10.1029/2005JD006548.
- 531 Collins, M., B. B. Booth, B. Bhaskaran, G. R. Harris, J. M. Murphy, D. M. H. Sexton,  
532 and M. J. Webb, 2011: Climate model errors, feedbacks and forcings: a comparison of  
533 perturbed physics and multi-model ensembles. *Clim. Dyn.*, **36**, 1737–1766, doi:10.1007/  
534 s00382-010-0808-0.
- 535 Ferranti, L. and P. Viterbo, 2006: The european summer of 2003: Sensitivity to soil water  
536 initial conditions. *J. Climate*, **19**, 3659–3680, doi:10.1175/JCLI3810.1.
- 537 Goddard, L., et al., 2013: A verification framework for interannual-to-decadal predictions  
538 experiments. *Climate Dynamics*, **40**, 245–272, doi:10.1007/s00382-012-1481-2.
- 539 Gordon, C., C. Cooper, C. A. Senior, H. Banks, J. M. Gregory, T. C. Johns, J. F. B. Mitchell,  
540 and R. A. Wood, 2000: The simulation of SST, sea ice extents and ocean heat transports  
541 in a version of the Hadley Centre coupled model without flux adjustments. *Climate Dyn.*,  
542 **16**, 147–168.
- 543 Gregory, J. M. and P. M. Forster, 2008: Transient climate response estimated from radiative  
544 forcing and observed temperature change. *Journal of Geophysical Research: Atmospheres*,  
545 **113**, D23 105, doi:10.1029/2008JD010405.
- 546 Hansen, J., R. Ruedy, M. Sato, and K. Lo, 2010: Global surface temperature change. *Rev.*  
547 *Geophys.*, **48**, RG4004, doi:10.1029/2010RG000345.
- 548 Harris, G. R., D. M. Sexton, B. B. Booth, M. Collins, and J. M. Murphy, 2013: Probabilistic  
549 projections of transient climate change. *Climate Dynamics*, **40**, 2937–2972, doi:10.1007/  
550 s00382-012-1647-y.
- 551 Hawkins, E. and R. Sutton, 2012: Time of emergence of climate signals. *Geophysical Research*  
552 *Letters*, **39**, L01 702, doi:10.1029/2011GL050087.
- 553 Ho, C. K., E. Hawkins, L. Shaffrey, J. Bröcker, L. Hermanson, J. M. Murphy, D. M. Smith,  
554 and R. Eade, 2013: Examining reliability of seasonal to decadal sea surface temperature  
555 forecasts: the role of ensemble dispersion. *Geophys. Res. Lett.*, **in press**, doi:10.1002/  
556 2013GL057630.
- 557 Kalnay, E., et al., 1996: The NCEP/NCAR 40-year reanalysis project. *BAMS*, **77**, 437–471.

- 558 Keenlyside, N. S., M. Latif, J. Jungclaus, L. Kornblueh, and E. Roeckner, 2008: Advancing  
559 decadal-scale climate prediction in the North Atlantic sector. *Nature*, **453**, 84–88, doi:  
560 10.1038/nature06921.
- 561 Knutti, R. and L. Tomassini, 2008: Constraints on the transient climate response from  
562 observed global temperature and ocean heat uptake. *Geophysical Research Letters*, **35**,  
563 L09 701, doi:10.1029/2007GL032904.
- 564 Masson, D. and R. Knutti, 2013: Predictor screening, calibration, and observational con-  
565 straints in climate model ensembles: an illustration using climate sensitivity. *J. Climate*,  
566 **26**, 887–898, doi:10.1175/JCLI-D-11-00540.1.
- 567 Meehl, G. A., L. Goddard, G. Boer, R. Burgman, and et al., 2013: Decadal climate predic-  
568 tion: An update from the trenches. *BAMS*, in press, doi:10.1175/BAMS-D-12-00241.1.
- 569 Meehl, G. A., et al., 2009: Decadal prediction: can it be skillful? *Bull. Amer. Met. Soc.*, **90**,  
570 1467–1485, doi:10.1175/2009BAMS2607.1.
- 571 Morice, C. P., J. J. Kennedy, N. A. Rayner, and P. D. Jones, 2012: Quantifying un-  
572 certainties in global and regional temperature change using an ensemble of observa-  
573 tional estimates: The HadCRUT4 data set. *J. Geophys. Res. Atmos.*, **117**, D08 101, doi:  
574 10.1029/2011JD017187.
- 575 Murphy, D. M., 2010: Constraining climate sensitivity with linear fits to outgoing radiation.  
576 *Geophysical Research Letters*, **37**, L09 704, doi:10.1029/2010GL042911.
- 577 Murphy, J. M., D. M. H. Sexton, D. N. Barnett, G. S. Jones, M. J. Webb, M. Collins, and  
578 D. A. Stainforth, 2004: Quantification of modelling uncertainties in a large ensemble of  
579 climate change simulations. *Nature*, **430**, 768–772, doi:10.1038/nature02771.
- 580 Oreskes, N., D. A. Stainforth, and L. A. Smith, 2010: Adaptation to global warming: Do  
581 climate models tell us what we need to know? *Philosophy of Science*, **77 (5)**, 1012–1028,  
582 doi:10.1086/657428.
- 583 Otto, A., F. Otto, O. Boucher, and et al., 2013: Energy budget constraints on climate  
584 response. *Nature Geosci.*, **6**, 415–416, doi:10.1038/ngeo1836.
- 585 Pohlmann, H., J. Jungclaus, A. Kohl, D. Stammer, and J. Marotzke, 2009: Initializing  
586 decadal climate predictions with the GECCO oceanic synthesis: Effects on the North  
587 Atlantic. *J. Climate*, **22**, 3926–3938, doi:10.1175/2009JCLI2535.1.
- 588 Randall, D. A., et al., 2007: *Climate Models and Their Evaluation. In: Climate Change*  
589 *2007: The Physical Science Basis*. Cambridge University Press, Cambridge, UK.

- 590 Raper, S. C. B., J. M. Gregory, and R. J. Stouffer, 2002: The role of climate sensitivity and  
591 ocean heat uptake on aogcm transient temperature response. *J. Climate*, **15**, 124–130,  
592 doi:10.1175/1520-0442(2002)015<0124:TROCSA>2.0.CO;2.
- 593 Robson, J. I., 2010: Understanding the performance of a decadal prediction system. Ph.D.  
594 thesis, University of Reading, UK.
- 595 Robson, J. I., R. T. Sutton, and D. M. Smith, 2012: Initialized decadal predictions of the  
596 rapid warming of the north atlantic ocean in the mid 1990s. *Geophysical Research Letters*,  
597 **39**, L19713, doi:10.1029/2012GL053370.
- 598 Smith, D. M., S. Cusack, A. W. Colman, C. K. Folland, G. R. Harris, and J. M. Murphy,  
599 2007: Improved surface temperature prediction for the coming decade from a global climate  
600 model. *Science*, **317**, 796–799, doi:10.1126/science.1139540.
- 601 Smith, D. M., R. Eade, N. J. Dunstone, D. Fereday, J. M. Murphy, H. Pohlmann, and  
602 A. Scaife, 2010: Skilful multi-year predictions of Atlantic hurricane frequency. *Nature*  
603 *Geosci.*, **3**, 846–849, doi:10.1038/ngeo1004.
- 604 Smith, D. M., et al., 2013: Real-time multi-model decadal climate predictions. *Climate*  
605 *Dynamics*, in press, doi:10.1007/s00382-012-1600-0.
- 606 Stocker, T. F., et al., (Eds.) , 2013: *IPCC, 2013: Summary for Policymakers. In: Cli-*  
607 *mate Change 2013: The Physical Science Basis. Contribution of Working Group I to the*  
608 *Fifth Assessment Report of the Intergovernmental Panel on Climate Change*. Cambridge  
609 University Press, Cambridge, UK.
- 610 Stott, P. A. and C. E. Forest, 2007: Review. Ensemble climate predictions using climate  
611 models and observational constraints. *Phil. Trans. R. Soc. A*, **365**, 2029–2052, doi:10.  
612 1098/rsta.2007.2075.
- 613 Taylor, K. E., R. J. Stouffer, and G. A. Meehl, 2012: An overview of CMIP5 and the exper-  
614 iment design. *Bull. Amer. Meteor. Soc.*, **93**, 485–498, doi:10.1175/BAMS-D-11-00094.1.
- 615 Tett, S., D. Rowlands, M. Mineter, and C. Cartis, 2013: Can top of atmosphere radiation  
616 measurements constrain climate predictions? part 2: Climate sensitivity. *J. Climate*, in  
617 press.
- 618 Toniazzo, T. and S. Woolnough, 2013: Development of warm SST errors in the south-  
619 ern tropical Atlantic in CMIP5 decadal hindcasts. *Clim. Dyn.*, in press, doi:10.1007/  
620 s00382-013-1691-2.

- 621 Uppala, S. M., et al., 2005: The ERA-40 re-analysis. *QJRMS*, **131**, 2961–3012, doi:10.1256/  
622 qj.04.176.
- 623 van Oldenborgh, G. J., F. J. Doblas-Reyes, B. Wouters, and W. Hazeleger, 2012: Decadal  
624 prediction skill in a multi-model ensemble. *Climate Dynamics*, **38**, 1263–1280, doi:10.  
625 1007/s00382-012-1313-4.
- 626 Yeager, S., A. Karspeck, G. Danabasoglu, J. Tribbia, and H. Teng, 2012: A decadal pre-  
627 diction case study: Late 20th century north atlantic ocean heat content. *J. Climate*, **25**,  
628 5173–5189.

629 **List of Tables**

630 1 The 5-95% ranges and medians (in brackets) of the original TCR (K) and the  
631 bias constrained values using a Monte-Carlo approach of linear fits to TCR  
632 against different observations. 22

TABLE 1. The 5-95% ranges and medians (in brackets) of the original TCR (K) and the bias constrained values using a Monte-Carlo approach of linear fits to TCR against different observations.

	TCR
Original	1.61-2.64 (2.17)
Constrained ranges	
ERA40	1.65-1.99 (1.82)
NCEP	1.59-1.91 (1.75)
GISS	1.61-1.93 (1.77)
HadCRUT4	1.45-1.83 (1.64)

## 633 List of Figures

- 634 1 A schematic illustrating the definition of  $B_{\text{obsvar}}$  (Eqn. 7) and consistent verifi-  
635 cation times (Section 2e). (a-c) Black lines show pseudo-observations, the red  
636 lines show pseudo-predictions (with noise) for three lead times ( $\tau$ ) as labelled,  
637 and the grey regions indicate the area integrated in the definition of  $B_{\text{obsvar}}$ .  
638 The blue bars indicate the range of times which are considered ‘consistent’,  
639 i.e. where all lead times can be simultaneously assessed. (d)  $B_{\text{obsvar}}$  for all  
640 verification times (black) and consistent verification times (blue). (e) Same  
641 as (d) for  $B'_{\text{obsvar}}$ . 27
- 642 2 (a) Example of a simple pseudo-prediction system, including observations  
643 (black), predictions (red), the true forced trend (solid blue) and estimated  
644 forced trend (dashed blue). (b) The bias tendency estimates for the predic-  
645 tions in (a), showing the true bias tendency (dark grey), the raw bias tendency  
646 estimate (solid blue), the bias tendency corrected using observed variability  
647 ( $B_{\text{obsvar}}$ ) for the cases when the forced trend is known (black) and unknown  
648 (dashed blue). (c) 10 examples of the bias estimates in (b) with different  
649 realisations of the observations. 28
- 650 3 The spread in 1000 realisations of the bias tendency estimates, an example of  
651 which is shown in Fig. 2, for the raw bias tendency (solid black) and corrected  
652 bias tendency (dashed black) at a lead time of 10 years. The magnitude of  
653 the true bias is shown in grey, indicating that, for this choice of toy model  
654 parameters, the bias could be detected with  $L \approx 16$  (20) hindcast start dates  
655 if the correction is made (or not). 29

- 656 4 (a) Example of a pseudo-prediction system with a lead time dependent bias,  
657 including observations (black), hindcasts (red), the true forced trend (solid  
658 blue) and estimated forced trend (dashed blue), including a mock volcanic  
659 eruption. (b) The bias tendency estimates for the predictions in (a), showing  
660 the true bias tendency (dark grey), true forcing bias tendency (light grey),  
661 the raw bias tendency estimates (blue), the bias tendency using consistent  
662 verification times (red) and the bias tendency estimates corrected using the  
663 consistent bias tendency (green). The dashed blue and green lines are cor-  
664 rected using  $B'_{\text{obsvar}}$ . 30
- 665 5 Spread in bias tendency estimates at a lead time of 10 years, as a function  
666 of the number of ensemble members considered, for (a) fully corrected bias  
667 estimate, (b) no observed variability correction, (c) no lead time dependent  
668 forcing bias correction, and (d) the raw bias. 31
- 669 6 An estimate of the start-time dependent forcing bias in the NoAssim predic-  
670 tion system (Smith et al. 2010). The left column shows the forcing estimates  
671 used in the transient integrations, the middle column shows the estimated  
672 forcing used in NoAssimPPE, and the right column shows the difference. The  
673 eruptions of Agung, El Chichon and Pinatubo are the main cause of the bias. 32
- 674 7 Toy model estimates for the error in true bias tendency estimates for the hind-  
675 cast setup of two operational prediction systems, namely NoAssim1 (Smith  
676 et al. 2007) and NoAssim PPE (Smith et al. 2010). NoAssim1 uses 20 years of  
677 hindcasts, with an effective ensemble size of 16 members (black line). NoAs-  
678 sim PPE uses 40 years of hindcasts with 9 different perturbed physics versions  
679 of the model, each with a single member. These can be considered as inde-  
680 pendent single member ensembles (solid grey) or as a 9-member ensemble  
681 (dashed grey). The horizontal error bars indicate the errors for 5 year mean  
682 predictions for NoAssimPPE (single members). 33

- 683 8 Bias tendency estimates (K) for global mean surface air temperature using  
684 NoAssim PPE. Different colours represent different observational datasets.  
685 (a) Raw bias. (b) Consistent verification times bias which is an estimate of  
686 the start-time dependent forcing bias. (c) Raw bias corrected by obsvar. (d)  
687 The true bias estimate, which is (c)-(b). The error ranges in (d) are derived  
688 from the toy model (Fig. 7) and are shown relative to the ERA-40 results. 34
- 689 9 Time evolution of ensemble mean (a) true bias tendency ( $\text{Wm}^{-2}$ ) in shortwave  
690 radiation at the top of atmosphere (TOA) of HadCM3 NoAssim PPE hind-  
691 casts for the period 1961-2001 against ERA-40 data set, (b) implied surface  
692 heat flux bias tendency ( $\text{Wm}^{-2}$ ) from integrated ocean heat content (OHC)  
693 bias for top 1500 metres against the Met Office ocean analysis and (c) esti-  
694 mated global mean error ( $\text{Wm}^{-2}$ ) associated with volcanic forcing in hindcasts. 35
- 695 10 The components of the total bias tendency for NoAssimPPE against Had-  
696 CRUT4 data. The total bias tendency (black) is dominated by the lead-time  
697 dependent forcing bias (green). The magnitude of the forcing bias is qualita-  
698 tively consistent with the magnitude of the forcing errors (Fig. 6). 36
- 699 11 Time evolutions of ensemble mean bias tendencies (K) for ocean temperature  
700 at 5m and ocean heat content (top 100m and top 500m) of HadCM3 NoAs-  
701 sim PPE hindcasts for the period 1961-2010 against Met Office ocean analysis  
702 data. (a) using all verification times (1961-2010), (b), using consistent verifi-  
703 cation times (1971-2001), (c) true bias tendency with linear trend removed in  
704 the analysis before calculating bias tendency associated with observed vari-  
705 ability. (d) Time evolution of ensemble mean true bias tendency (K) as a  
706 function of depth for global ocean temperature for HadCM3 NoAssim PPE  
707 hindcasts for the period 1961-2010 against the Met Office ocean analysis. 37

- 708 12 Relationships between global mean SAT true bias tendencies (K) (against  
709 HadCRUT4 data) and global mean OHC (top 1000m) bias tendencies (against  
710 the Met Office ocean analysis) for 9 PPE model versions. (a) average for lead  
711 years 1-5, and (b) average for lead years 6-10. 38
- 712 13 Relationships between the lead years 6-10 averaged global mean SAT true  
713 bias tendencies (K) against HadCRUT4 data for each version of PPE hind-  
714 casts for (a) TCR and (b) non-GHG aerosol forcing trend, using 9 PPE model  
715 versions. The error bars for bias tendency are based on the toy model (Fig. 7).  
716 Grey lines are example linear fits to TCR and to the non-GHG aerosol forc-  
717 ing trend using a Monte-Carlo approach, and the red lines are the best fit.  
718 The constrained ranges of TCR and the non-GHG aerosol forcing trend are  
719 shown as black bars assuming a true bias tendency error of 0.016K (solid)  
720 and 0.032K (dashed). Other ranges for TCR (Stocker et al. 2013, Gregory and  
721 Forster 2008 - denoted GF08) ranges are also given. (c,d) estimated proba-  
722 bility distribution functions (PDFs) of unconstrained (blue) and constrained  
723 (full black and dotted black) TCR and non-GHG aerosol forcing trends. The  
724 dashed black lines indicate the PDF for doubled uncertainties in the true bias  
725 tendency. 39

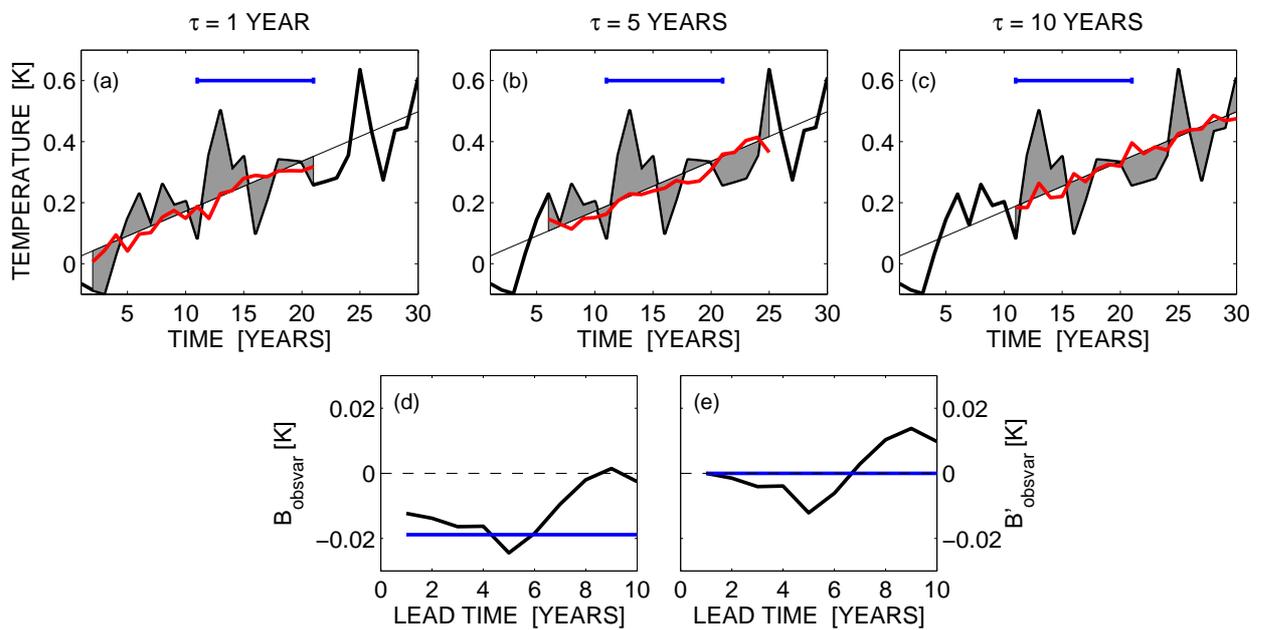


FIG. 1. A schematic illustrating the definition of  $B_{\text{obsvar}}$  (Eqn. 7) and consistent verification times (Section 2e). (a-c) Black lines show pseudo-observations, the red lines show pseudo-predictions (with noise) for three lead times ( $\tau$ ) as labelled, and the grey regions indicate the area integrated in the definition of  $B_{\text{obsvar}}$ . The blue bars indicate the range of times which are considered ‘consistent’, i.e. where all lead times can be simultaneously assessed. (d)  $B_{\text{obsvar}}$  for all verification times (black) and consistent verification times (blue). (e) Same as (d) for  $B'_{\text{obsvar}}$ .

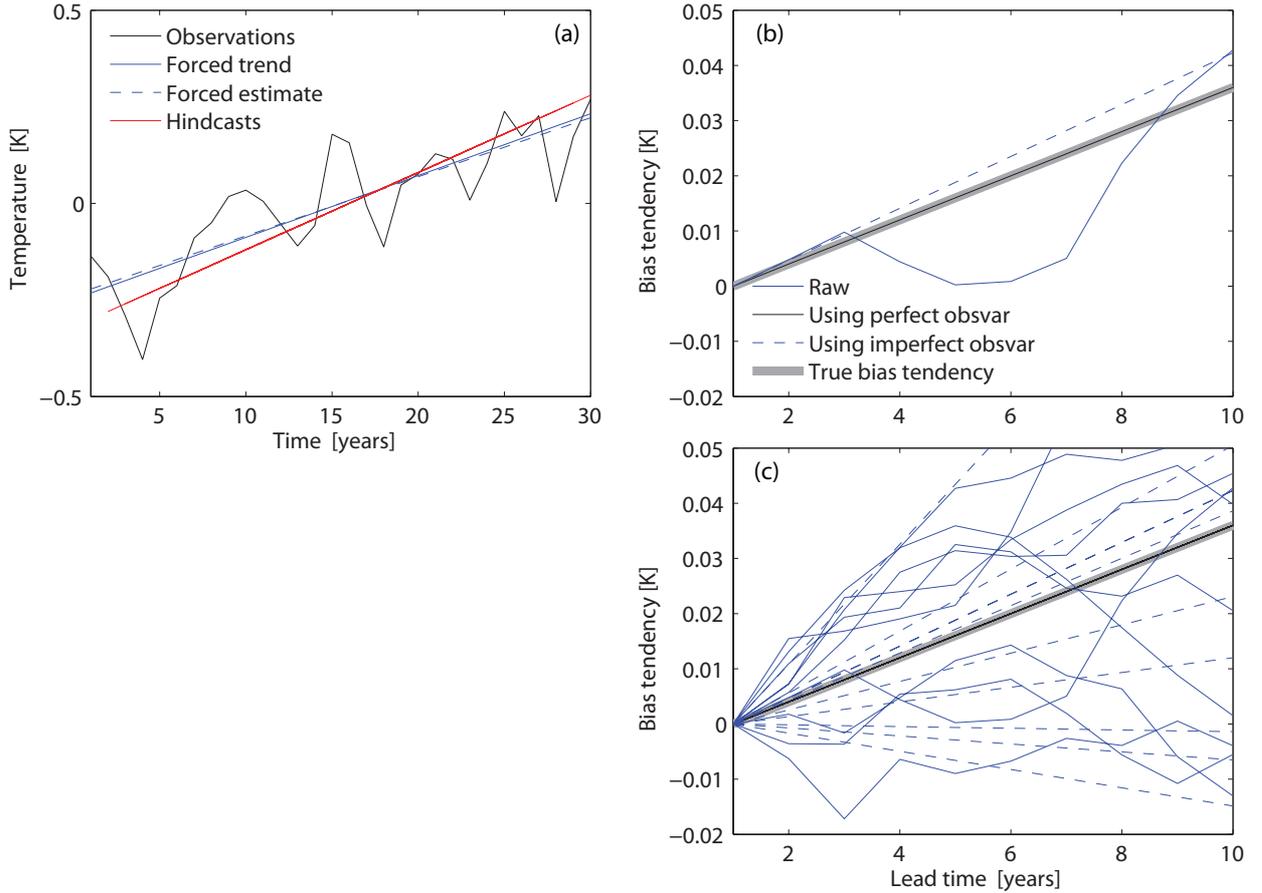


FIG. 2. (a) Example of a simple pseudo-prediction system, including observations (black), predictions (red), the true forced trend (solid blue) and estimated forced trend (dashed blue). (b) The bias tendency estimates for the predictions in (a), showing the true bias tendency (dark grey), the raw bias tendency estimate (solid blue), the bias tendency corrected using observed variability ( $B_{\text{obsvar}}$ ) for the cases when the forced trend is known (black) and unknown (dashed blue). (c) 10 examples of the bias estimates in (b) with different realisations of the observations.

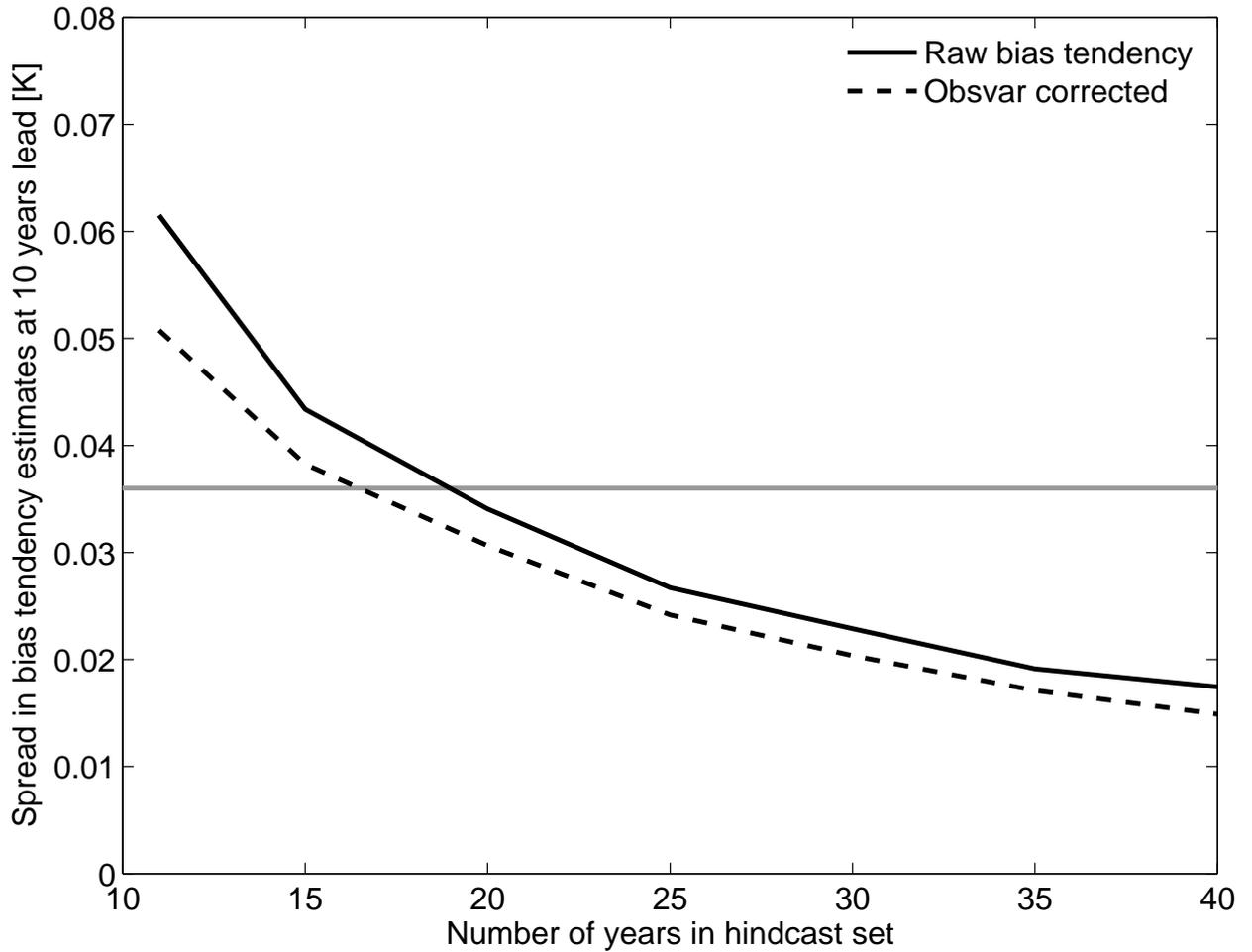


FIG. 3. The spread in 1000 realisations of the bias tendency estimates, an example of which is shown in Fig. 2, for the raw bias tendency (solid black) and corrected bias tendency (dashed black) at a lead time of 10 years. The magnitude of the true bias is shown in grey, indicating that, for this choice of toy model parameters, the bias could be detected with  $L \approx 16$  (20) hindcast start dates if the correction is made (or not).

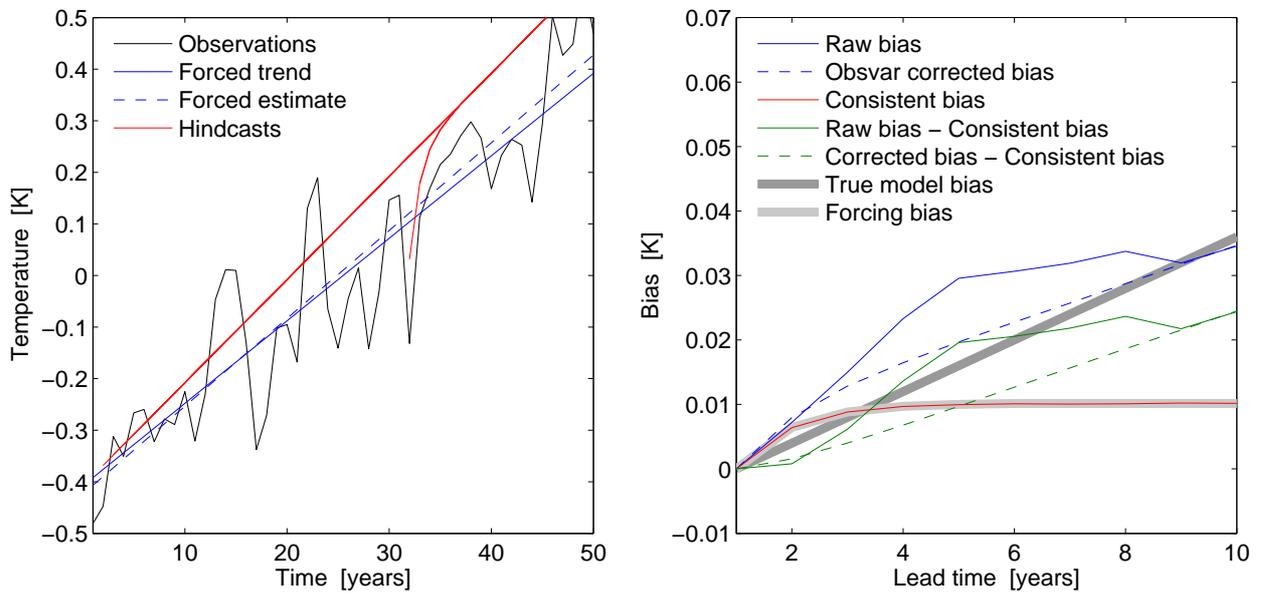


FIG. 4. (a) Example of a pseudo-prediction system with a lead time dependent bias, including observations (black), hindcasts (red), the true forced trend (solid blue) and estimated forced trend (dashed blue), including a mock volcanic eruption. (b) The bias tendency estimates for the predictions in (a), showing the true bias tendency (dark grey), true forcing bias tendency (light grey), the raw bias tendency estimates (blue), the bias tendency using consistent verification times (red) and the bias tendency estimates corrected using the consistent bias tendency (green). The dashed blue and green lines are corrected using  $B'_{\text{obsvar}}$ .

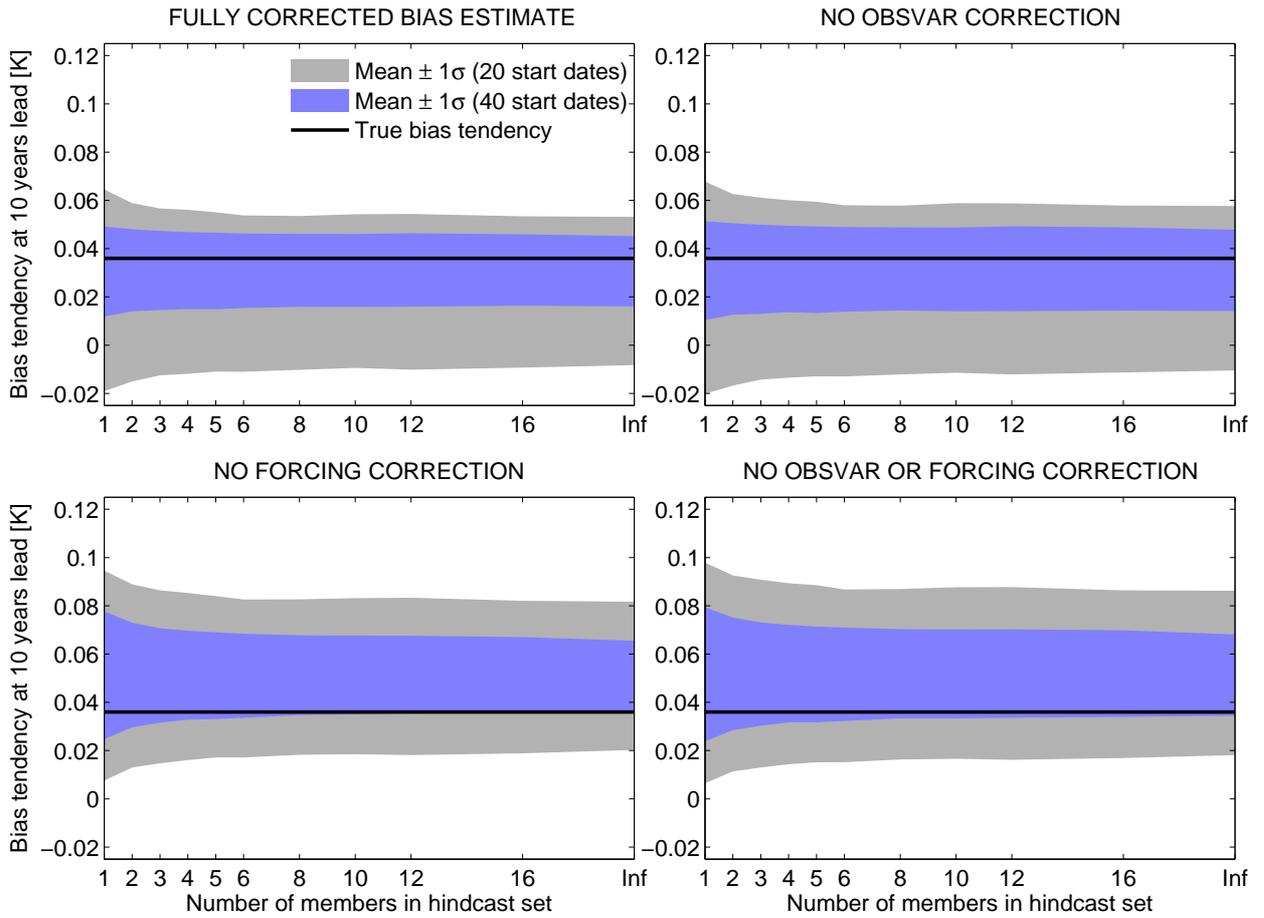


FIG. 5. Spread in bias tendency estimates at a lead time of 10 years, as a function of the number of ensemble members considered, for (a) fully corrected bias estimate, (b) no observed variability correction, (c) no lead time dependent forcing bias correction, and (d) the raw bias.

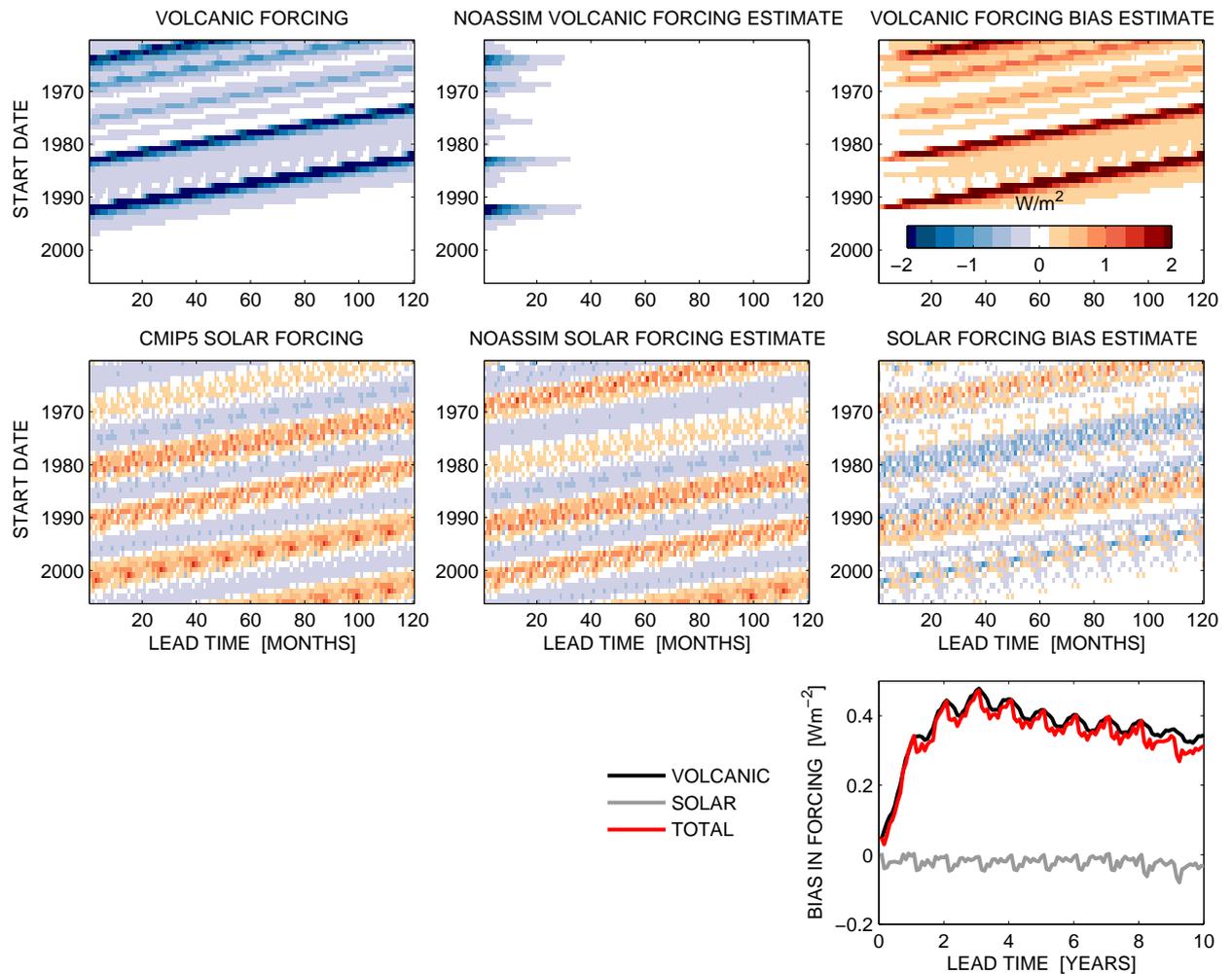


FIG. 6. An estimate of the start-time dependent forcing bias in the NoAssim prediction system (Smith et al. 2010). The left column shows the forcing estimates used in the transient integrations, the middle column shows the estimated forcing used in NoAssimPPE, and the right column shows the difference. The eruptions of Agung, El Chichon and Pinatubo are the main cause of the bias.

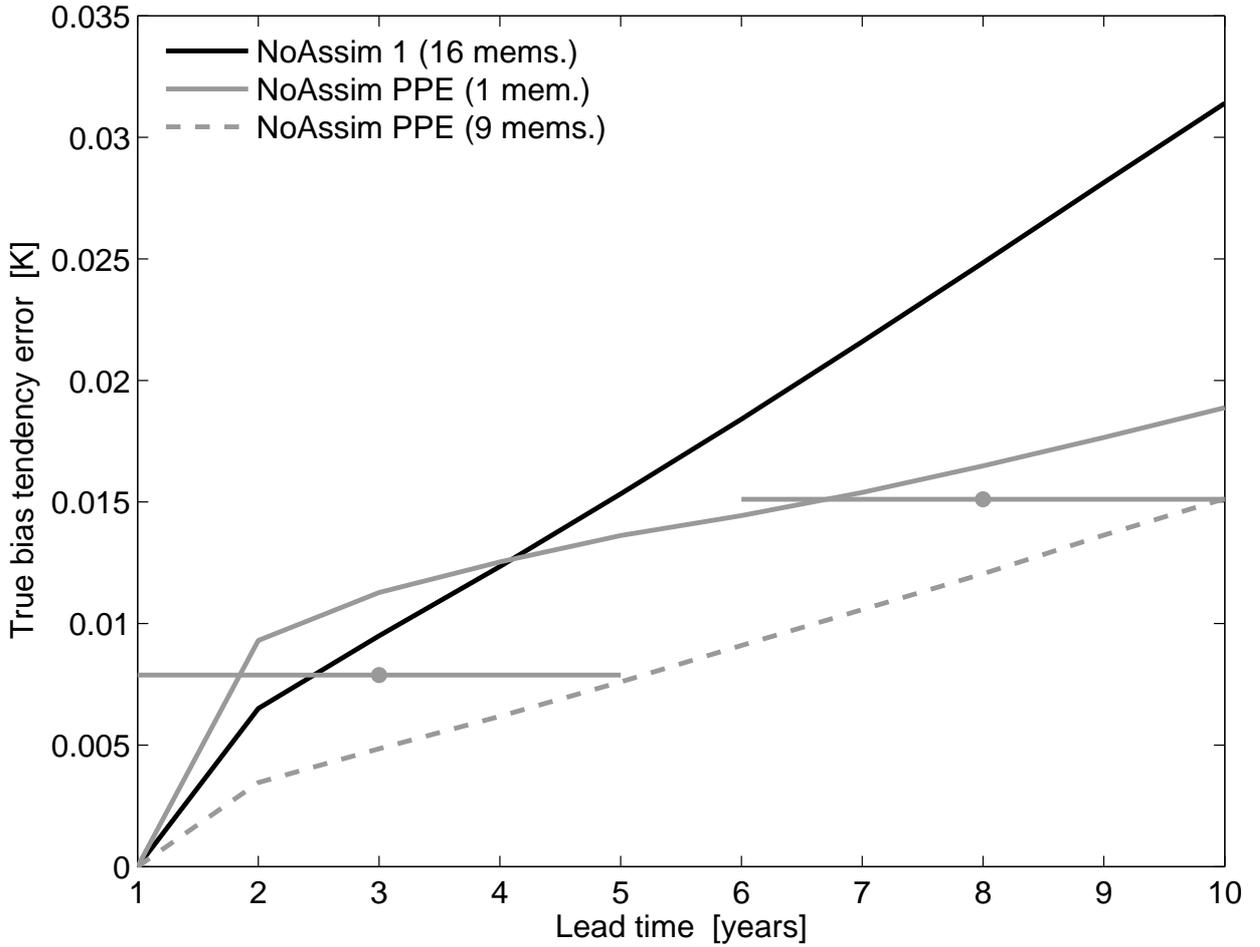


FIG. 7. Toy model estimates for the error in true bias tendency estimates for the hindcast setup of two operational prediction systems, namely NoAssim1 (Smith et al. 2007) and NoAssim PPE (Smith et al. 2010). NoAssim1 uses 20 years of hindcasts, with an effective ensemble size of 16 members (black line). NoAssim PPE uses 40 years of hindcasts with 9 different perturbed physics versions of the model, each with a single member. These can be considered as independent single member ensembles (solid grey) or as a 9-member ensemble (dashed grey). The horizontal error bars indicate the errors for 5 year mean predictions for NoAssimPPE (single members).

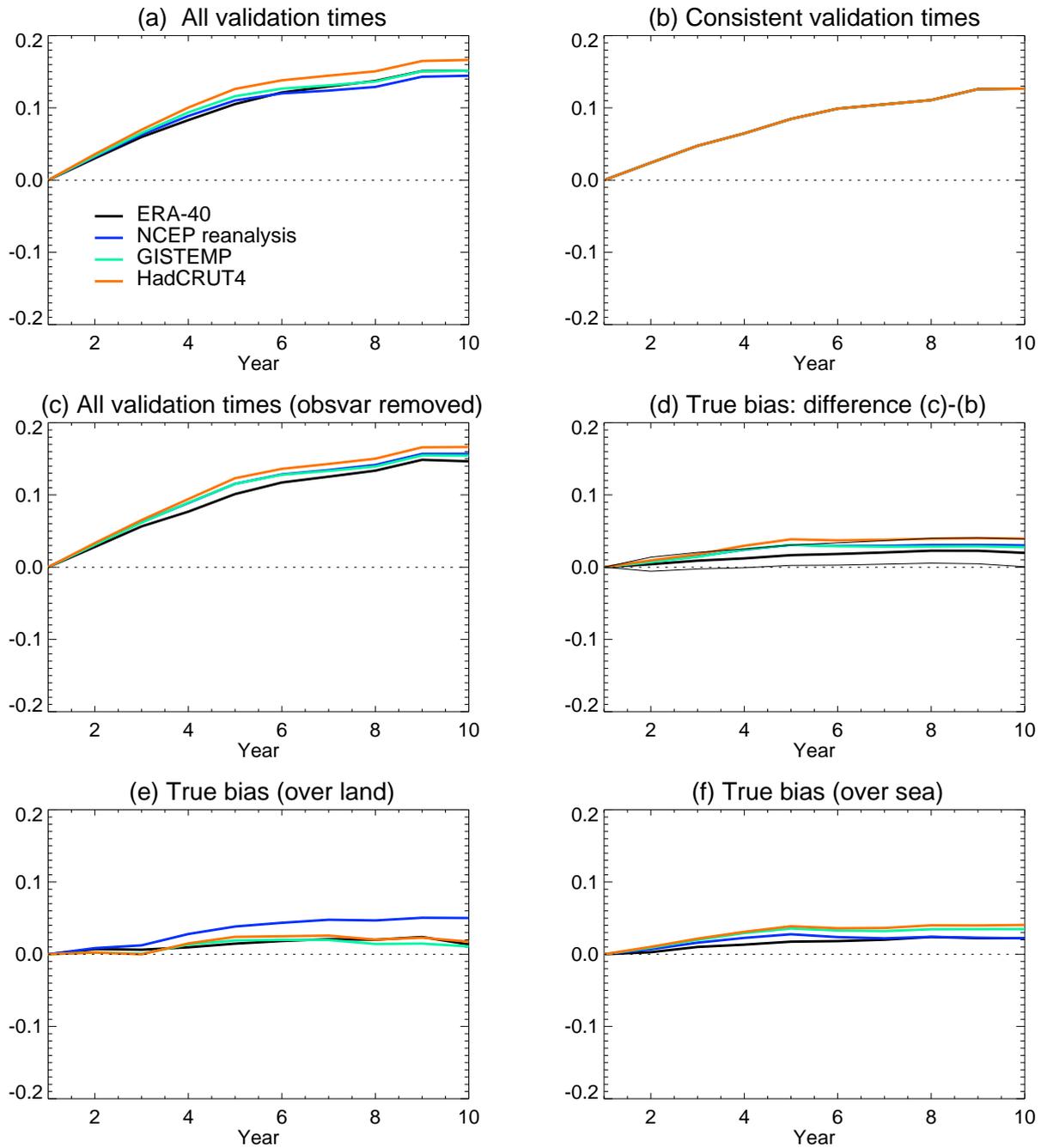
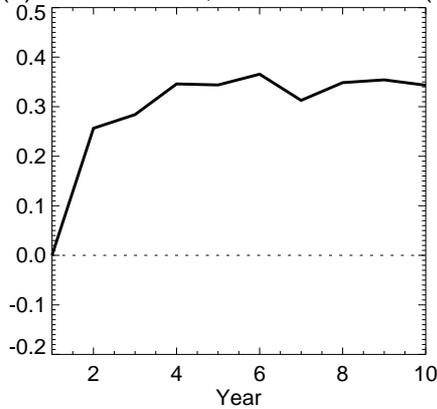
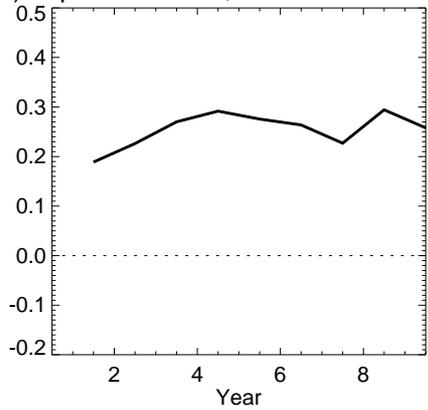


FIG. 8. Bias tendency estimates (K) for global mean surface air temperature using NoAssim PPE. Different colours represent different observational datasets. (a) Raw bias. (b) Consistent verification times bias which is an estimate of the start-time dependent forcing bias. (c) Raw bias corrected by obsvar. (d) The true bias estimate, which is (c)-(b). The error ranges in (d) are derived from the toy model (Fig. 7) and are shown relative to the ERA-40 results.

(a) TOA SW bias, consistent times (sea)



(b) Implied flux bias, consistent times (sea)



(c) Estimated volcanic forcing error

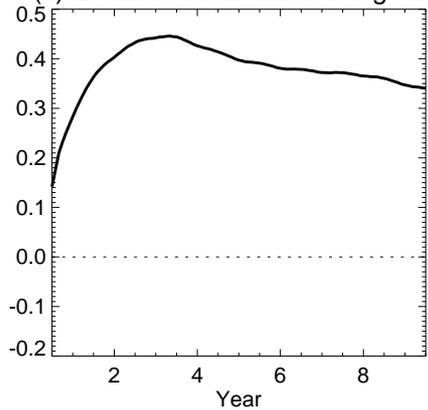


FIG. 9. Time evolution of ensemble mean (a) true bias tendency ( $\text{Wm}^{-2}$ ) in shortwave radiation at the top of atmosphere (TOA) of HadCM3 NoAssim PPE hindcasts for the period 1961-2001 against ERA-40 data set, (b) implied surface heat flux bias tendency ( $\text{Wm}^{-2}$ ) from integrated ocean heat content (OHC) bias for top 1500 metres against the Met Office ocean analysis and (c) estimated global mean error ( $\text{Wm}^{-2}$ ) associated with volcanic forcing in hindcasts.

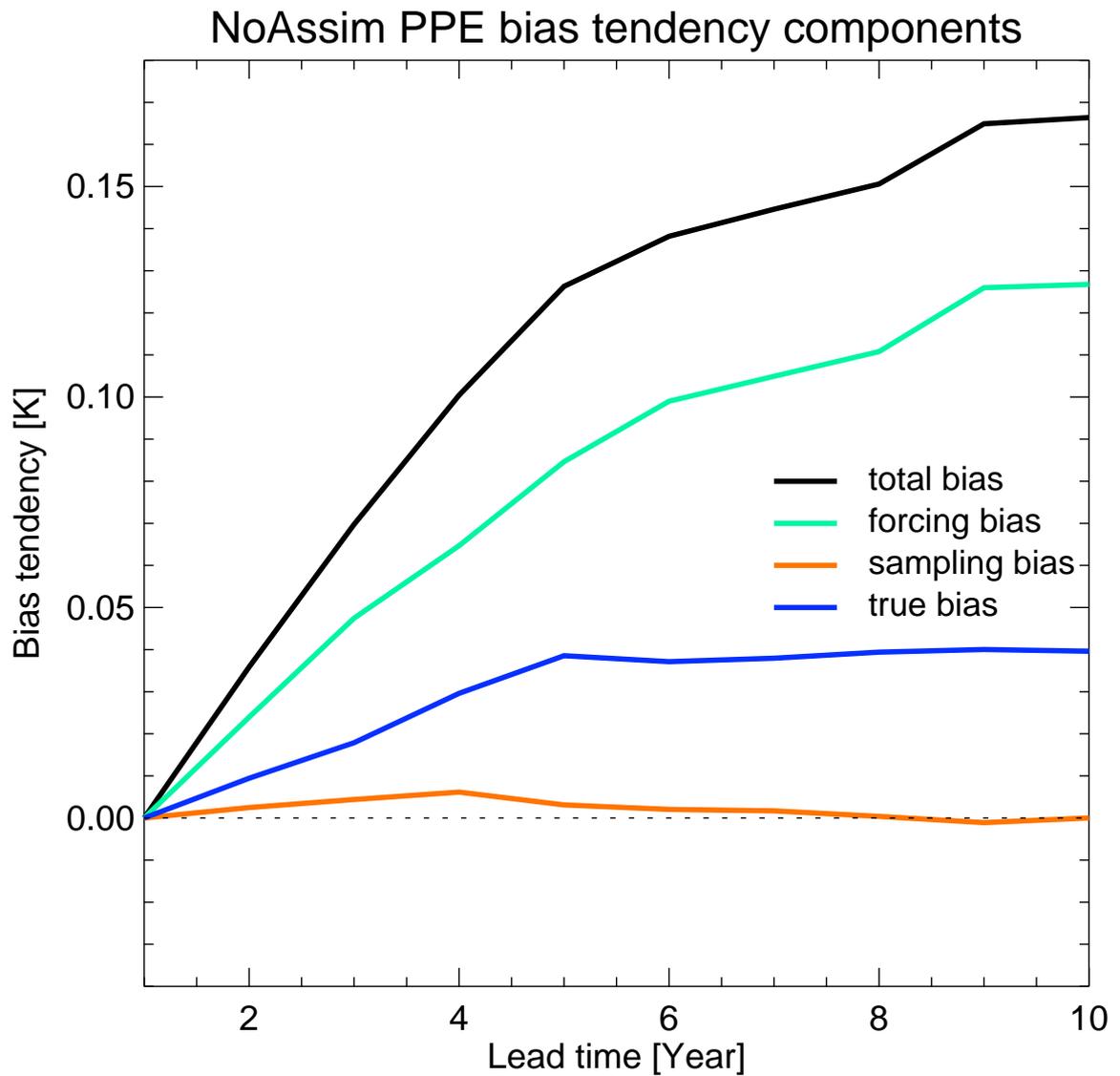


FIG. 10. The components of the total bias tendency for NoAssimPPE against HadCRUT4 data. The total bias tendency (black) is dominated by the lead-time dependent forcing bias (green). The magnitude of the forcing bias is qualitatively consistent with the magnitude of the forcing errors (Fig. 6).

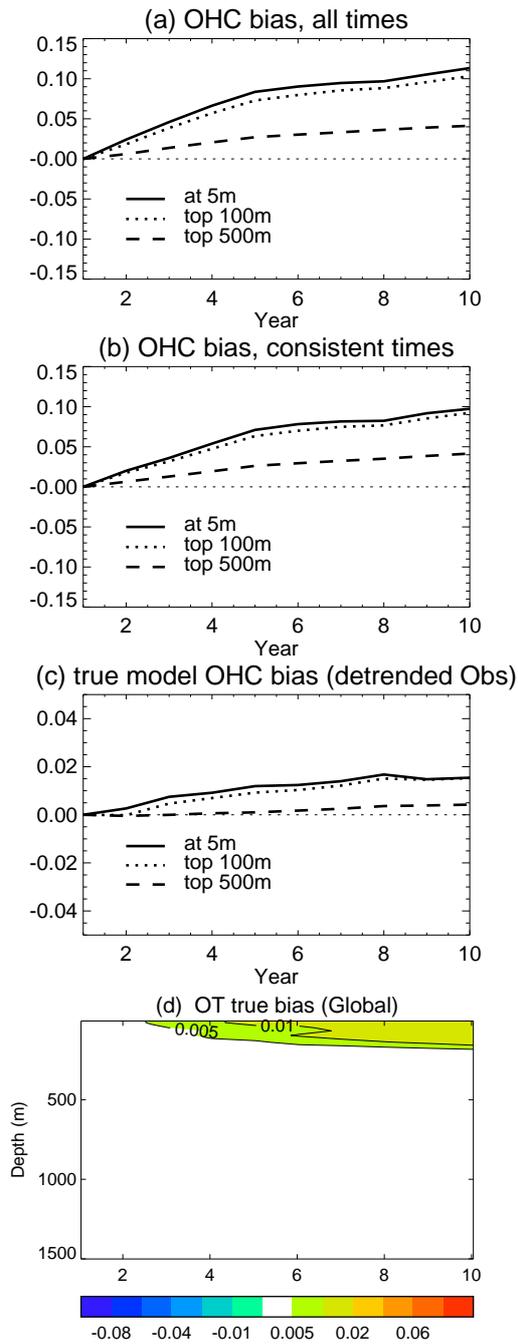


FIG. 11. Time evolutions of ensemble mean bias tendencies (K) for ocean temperature at 5m and ocean heat content (top 100m and top 500m) of HadCM3 NoAssim PPE hindcasts for the period 1961-2010 against Met Office ocean analysis data. (a) using all verification times (1961-2010), (b), using consistent verification times (1971-2001), (c) true bias tendency with linear trend removed in the analysis before calculating bias tendency associated with observed variability. (d) Time evolution of ensemble mean true bias tendency (K) as a function of depth for global ocean temperature for HadCM3 NoAssim PPE hindcasts for the period 1961-2010 against the Met Office ocean analysis.

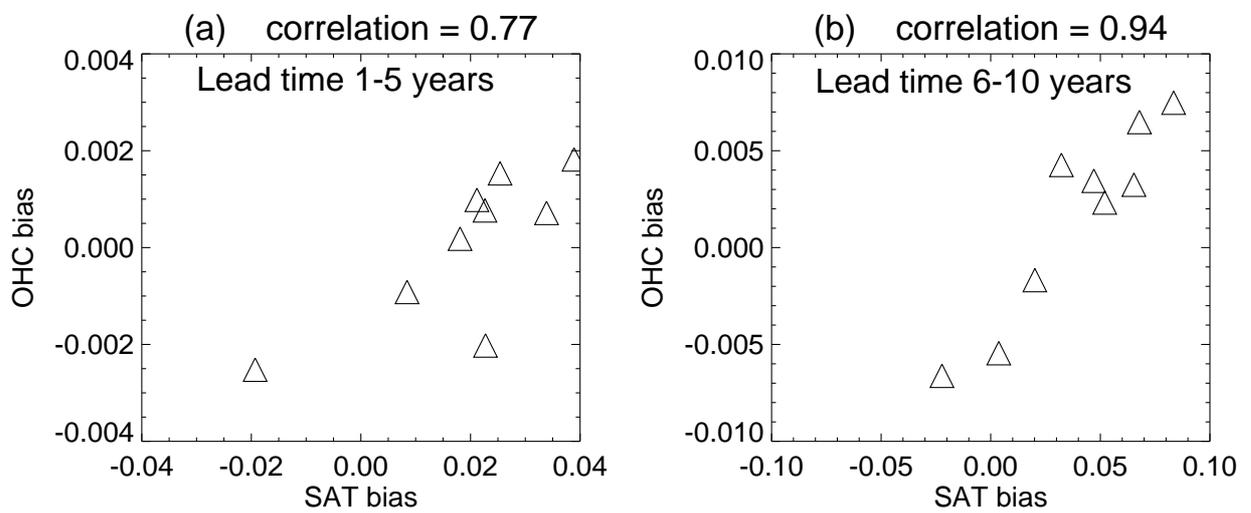


FIG. 12. Relationships between global mean SAT true bias tendencies (K) (against Had-CRUT4 data) and global mean OHC (top 1000m) bias tendencies (against the Met Office ocean analysis) for 9 PPE model versions. (a) average for lead years 1-5, and (b) average for lead years 6-10.

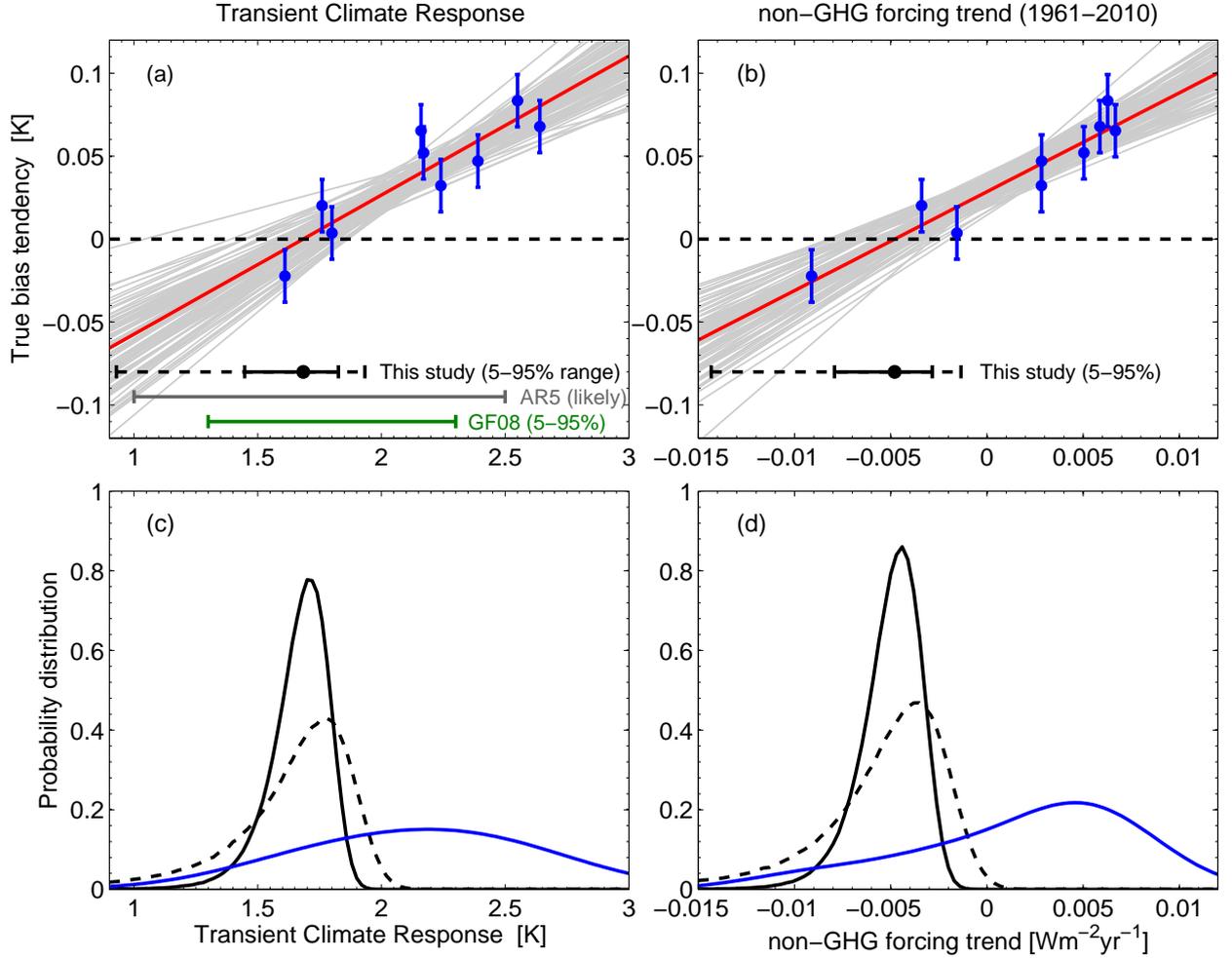


FIG. 13. Relationships between the lead years 6-10 averaged global mean SAT true bias tendencies (K) against HadCRUT4 data for each version of PPE hindcasts for (a) TCR and (b) non-GHG aerosol forcing trend, using 9 PPE model versions. The error bars for bias tendency are based on the toy model (Fig. 7). Grey lines are example linear fits to TCR and to the non-GHG aerosol forcing trend using a Monte-Carlo approach, and the red lines are the best fit. The constrained ranges of TCR and the non-GHG aerosol forcing trend are shown as black bars assuming a true bias tendency error of 0.016K (solid) and 0.032K (dashed). Other ranges for TCR (Stocker et al. 2013, Gregory and Forster 2008 - denoted GF08) ranges are also given. (c,d) estimated probability distribution functions (PDFs) of unconstrained (blue) and constrained (full black and dotted black) TCR and non-GHG aerosol forcing trends. The dashed black lines indicate the PDF for doubled uncertainties in the true bias tendency.