

Errors in simple climate model emulations of past and future global temperature change

L. S. Jackson¹, A. C. Maycock¹, T. Andrews², H.-B. Fredriksen³, C. J. Smith⁴, and P. M. Forster⁵

¹ School of Earth and Environment, University of Leeds, Leeds, UK.

² Met Office Hadley Centre, Exeter, UK.

³ Department of Physics and Technology, UiT the Arctic University of Norway, Tromsø, Norway.

⁴ International Institute for Applied Systems Analysis, Laxenburg, Austria.

⁵ Priestley International Centre for Climate, University of Leeds, Leeds, UK.

Corresponding author: L. S. Jackson (l.s.jackson@leeds.ac.uk)

Key Points:

- Emulators of global surface temperature calibrated to individual climate models can generate large errors in past and future predictions.
- Emulation errors are not systematically related to emulator parameters and vary between climate models meaning they cannot be predicted.
- Rigorous out-of-sample evaluation is necessary to characterize emulator performance.

18 Abstract

19 Climate model emulators are widely used to generate temperature projections for climate
20 scenarios, including in the recent IPCC Sixth Assessment Report. Here we evaluate the
21 performance of a two-layer energy balance model in emulating historical and future temperature
22 projections from CMIP6 models. We find that emulation errors can be large ($>0.5^{\circ}\text{C}$ for SSP2-
23 4.5) and differ markedly between climate models, forcing scenarios and time periods. Errors
24 arise in emulating the near-surface temperature response to both greenhouse gas and aerosol
25 forcing; in some periods the errors due to these forcings oppose one another, giving the spurious
26 impression of better emulator performance. Climate feedbacks are assumed constant in the
27 emulator, introducing time-varying or state dependent feedbacks may reduce prediction errors.
28 Close emulations can be produced for a given period but, crucially, this does not guarantee
29 reliable emulations of other scenarios and periods. Therefore, rigorous out-of-sample evaluation
30 is necessary to characterize emulator performance.

31 Plain Language Summary

32 Complex climate models are state-of-the-art tools used to produce projections of future
33 climate but they are expensive and take a long time to run. Climate model emulators are simple
34 statistical or physically based models that can aim to reproduce the response of complex climate
35 models to a prescribed climate change scenario at lower cost and more quickly. In this study, we
36 use a climate model emulator to reproduce simulations of twentieth and twenty-first century
37 temperatures for eight complex climate models. We show that close emulations can be produced
38 for pre-defined climate scenarios and time periods. Close emulations are not guaranteed,
39 however, when the emulator is used for other climate scenarios or other periods. This is
40 important because climate model emulators are frequently used to produce projections that are
41 not available from complex climate models. Evaluation of climate model emulators and
42 characterization of their uncertainties, therefore, should use data not used in the calibration of the
43 emulator.

44 **1 Introduction**

45 Climate model emulators are simplified physical or statistical models that are
46 computationally efficient. Climate model emulators played a central role in producing future
47 global near-surface temperature projections for Working Group I (Forster et al., 2021; Lee et al.,
48 2021) and Working Group III (Riahi et al., 2022) of the Sixth Assessment Report of the
49 Intergovernmental Panel on Climate Change (IPCC AR6). The IPCC AR6 used climate model
50 emulators to supplement simulations from coupled atmosphere-ocean general circulation models
51 (AOGCMs) extending available simulations further into the future and projecting future climate
52 scenarios not available from AOGCMs. It is important, therefore, that the simplifying
53 assumptions used by emulators are rigorously tested so the robustness of their performance is
54 understood.

55 Physically based climate model emulators, such as energy balance models (EBMs), use
56 bulk physical relationships to emulate the large-scale behavior of Earth’s climate system. For
57 example, EBMs were used by Colman and Soldatenko (2020) to investigate links between
58 climate variability and climate sensitivity and, by Modak and Mauritsen (2021) to investigate the
59 probability of occurrence of the 2000-2012 global warming hiatus.

60 Two-layer EBMs produce close emulations of idealized abrupt-4xCO₂ and 1pctCO₂
61 simulations from AOGCMs (e.g., “EBM- ϵ ” in Geoffroy et al. 2013b; “held-two-layer-uom” in
62 Nicholls et al., 2020). Differences between emulations and AOGCM projections are generally
63 greatest at times of pronounced change in the rate of temperature increase. Such changes are
64 associated with time-varying feedbacks (Senior and Mitchell, 2000; Winton et al., 2010; Armour
65 et al., 2013; Dong et al., 2020; Dunne et al., 2020; Rugenstein et al., 2020; Dong et al., 2021)
66 which are caused by evolving spatial pattern effects in surface temperature (Stevens et al., 2016;
67 Andrews et al., 2015; Rugenstein et al., 2016; Dong et al., 2021) and non-linear state
68 dependences in climate feedbacks (Good et al., 2015; Rohrschneider et al., 2019; Bloch-Johnson
69 et al., 2021). EBMs have been enhanced to capture time-varying feedbacks: the Geoffroy et al.
70 (2013b) EBM includes an efficacy parameter for deep ocean heat uptake and the “held-two-
71 layer-uom” EBM also includes a state dependent feedback parameter (Rohrschneider et al.,
72 2019; Nicholls et al., 2020). These paradigms, however, do not precisely capture the feedback
73 changes in AOGCMs and contribute to structural error which is irreducible unless the EBM

74 structure is enhanced (e.g., extending a two-layer EBM to three or more layers (Cummins et al.,
75 2020)).

76 Assessments of emulator performance are more trustworthy when projections are
77 validated using data different from those used to calibrate the emulator parameters (out-of-
78 sample validation). EBM parameters are frequently calibrated using idealized step-forcing
79 experiments (e.g., abrupt-4xCO₂) with the parameters estimated using analytical methods
80 (Geoffroy et al., 2013a) or statistical methods (e.g., Cummins et al., 2020). The Coupled Model
81 Intercomparison Project Phase 6 (CMIP6) (Eyring et al., 2016) historical and future shared
82 socio-economic pathway (SSP) projections for AOGCMs, therefore, are well suited for assessing
83 EBM emulator performance. They can be used to produce out-of-sample assessments using
84 realistic climate scenarios. Although climate model emulators have been evaluated (e.g.,
85 Nicholls et al., 2020; Nicholls et al., 2021), it is not known how well emulators perform for the
86 latest CMIP6 AOGCMs using realistic, out-of-sample climate projections and latest assessments
87 of effective radiative forcing (ERF). Furthermore, the contribution of irreducible structural errors
88 to total prediction error remains poorly understood.

89 In this study, we evaluate the performance of a two-layer energy balance model (EBM2)
90 (Held et al., 2010; Geoffroy et al., 2013a, b) for emulating CMIP6 historical and future
91 temperature trends using different EBM calibrations. We calibrate the EBM2 parameters for
92 specific periods and ERFs, and evaluate the temperature projections for subsequent periods and
93 alternative ERF scenarios. EBM2 is compared against a step-response emulator and a three-
94 layer EBM.

95 2 Methods and data

96 2.1 Step-response emulator

97 We use a step-response emulator (Good et al., 2011) to provide a comparator of EBM
98 emulator performance for temperature projections. The step-response function for each AOGCM
99 was derived by dividing the projected temperature changes from a single realization of a CMIP6
100 abrupt-4xCO₂ simulation by the radiative forcing for 4xCO₂ (Smith et al., 2020). The step-
101 response function was smoothed using cubic splines, and linear regression (years 121-150) was
102 used for extrapolation beyond the 150 years of the abrupt-4xCO₂ simulations. Temperature

103 projections from the step-response emulator were produced by convolution of annual changes in
 104 ERF and the step-response functions.

105 2.2 Two-layer EBM emulator (EBM2)

106 In EBM2 (Held et al., 2010; Geoffroy et al., 2013a) the upper layer represents the Earth's
 107 atmosphere, land surface and ocean mixed layer, and the lower layer represents the deep ocean.

108 The rate of temperature change in each EBM2 layer is determined from:

$$109 \quad C_1 \frac{dT_1}{dt} = F + \lambda T_1 - \varepsilon \gamma (T_1 - T_0) \quad (1)$$

$$110 \quad C_0 \frac{dT_0}{dt} = \gamma (T_1 - T_0) \quad (2)$$

111 Where C represents heat capacity, T temperature, F ERF, λ the climate feedback
 112 parameter and γ the heat transfer coefficient between the upper layer (layer 1) and the lower layer
 113 (layer 0). We follow the formulation of Geoffroy et al. (2013b) which includes an efficacy
 114 parameter for deep ocean heat uptake (ε) to account for the forced pattern effect in surface
 115 temperature (Stevens et al., 2016). As is commonplace (Geoffroy et al., 2013a, b; Gregory et al.,
 116 2015; Cummins et al., 2020), the EBM2 parameters were calibrated for each AOGCM using a
 117 single realization of a CMIP6 abrupt-4xCO2 simulation. Radiative forcing for 4xCO2 was taken
 118 from Smith et al. (2020). See Tables S1 and S2 for further details.

119 2.3 Calibration of EBM2 using linear optimization

120 As an alternative to calibration using the abrupt-4xCO2 experiment, we use linear
 121 optimization (the L-BFGS-B algorithm in `scipy.optimize.minimize` v1.6.2) to optimize the λ and
 122 ε parameters by minimizing the root mean square error (RMSE) of the emulated temperatures
 123 compared to the AOGCM over a defined time period (e.g., historical) (Table S3, S4). Lower
 124 bounds of $-0.5 \text{ W m}^{-2} \text{ K}^{-1}$ and 0.5 were imposed for λ and ε respectively, and upper bounds of -
 125 $2.0 \text{ W m}^{-2} \text{ K}^{-1}$ and 2.0 respectively. These bounds are broadly based on the range of parameter
 126 values from the abrupt-4xCO2 calibration. The temperature projections are less sensitive to
 127 changes in the other EBM2 parameters (i.e., C_0 , C_1 , and γ), so these parameters are unchanged
 128 from their abrupt-4xCO2 calibrations. We also applied the linear optimization methodology to
 129 the abrupt-4xCO2 simulations, which produced very similar parameter values to the Geoffroy et
 130 al. (2013b) methodology used in the abrupt-4xCO2 calibration.

131 2.4 Three-layer EBM

132 We use a three-layer EBM (EBM3) (Cummins et al., 2020) as a second comparator for
133 EBM2 performance. We follow the method of Cummins et al. (2020) to calibrate the EBM3
134 parameters (including ERF for 4xCO₂) using a single realization of a CMIP6 abrupt-4xCO₂
135 simulation.

136 2.5 Data

137 We use projections of global annual mean near-surface temperature and radiative fluxes
138 at the top of atmosphere (TOA) from the CMIP6 archive. We emulate temperatures for eight
139 AOGCMs selected because data was available for the CMIP6 experiments of interest. For
140 projections of recent and future climate change, the Historical and SSP experiments were used.
141 Projections of temperature change attributed to specific sources of ERF are taken from the
142 Detection and Attribution Model Intercomparison Project (DAMIP) experiments (Gillett et al.,
143 2016). The emulations are driven by time series of total annual ERF; estimates of ERF are taken
144 from the Radiative Forcing Model Intercomparison Project (RFMIP) experiments (Pincus et al.,
145 2016; Smith et al., 2021). The ERF for GFDL-CM4 was used for GFDL-ESM4 (RFMIP ERF
146 being unavailable for GFDL-ESM4). Following Forster et al. (2013), unforced drift is removed
147 from the AOGCM projections using the preindustrial control experiment.

148 3 Results

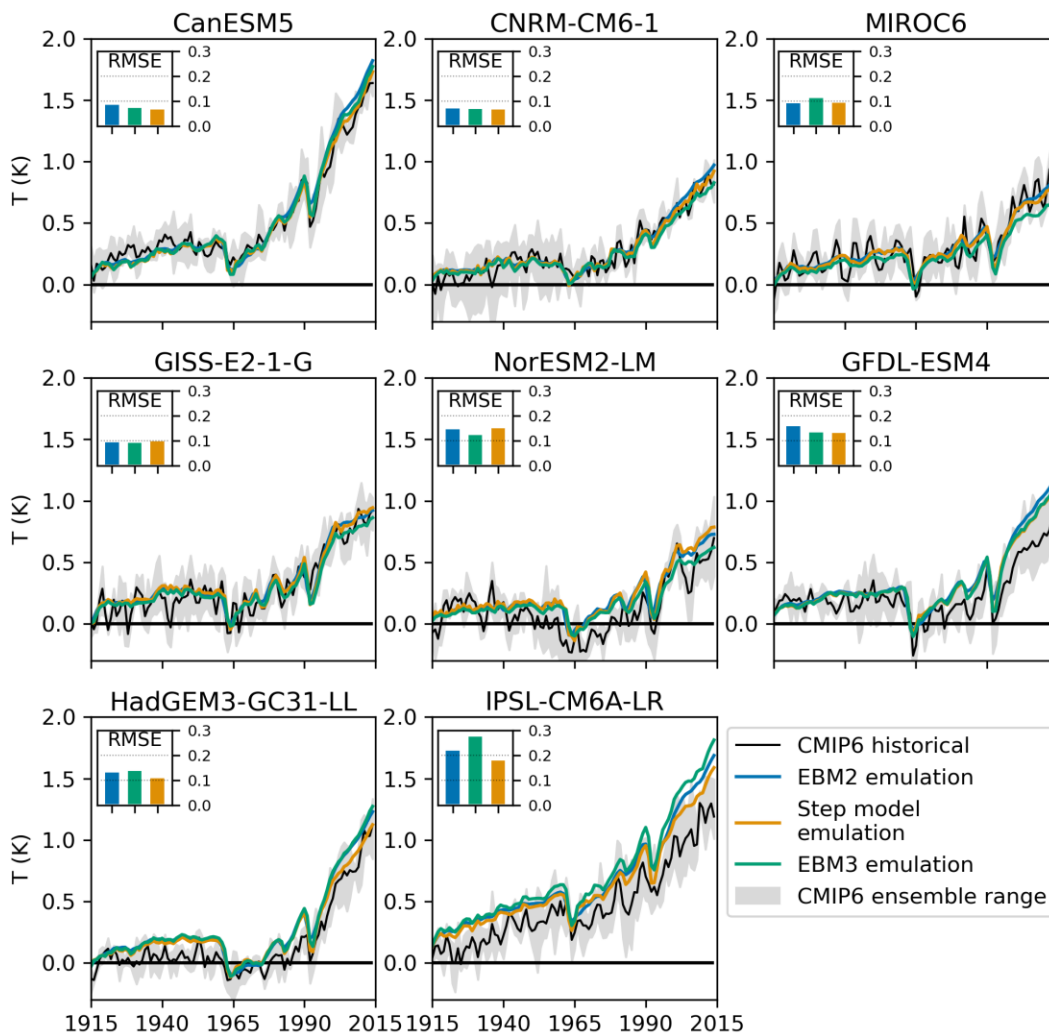
149 3.1 Historical period using the abrupt-4xCO₂ calibration

150 EBM2 captures the increasing temperature trend during the twentieth century and
151 distinguishes between high and low climate sensitivity AOGCMs (Figure 1). In all EBM2
152 emulations, a proportion of the RMSE (~ 0.07 K) arises from interannual variations in the
153 AOGCM ensemble means that is not captured in the emulations (there are up to three members
154 in each AOGCM historical ensemble). The performance of EBM2, however, varies substantially
155 between AOGCMs. The emulation errors are not strongly correlated with parameter values
156 though there is a weak correlation between smaller RMSEs and large relative deep ocean heat
157 capacity (i.e., C_o/C_l) (Figure S1). The sensitivity of emulation errors to changes in λ and ϵ varies
158 between AOGCMs (Figure S2). There are instances of both large and small RMSE emulations
159 for both high and low climate sensitivity AOGCMs. For AOGCMs where there are substantial

160 differences between the emulations and the AOGCM projections, the differences occur over
161 different time periods. Differences are large for 1925-1950 (HadGEM3-GC31-LL), for 1950-
162 1975 (NorESM2-LM) and for 2000-2015 (HadGEM3-GC31-LL, IPSL-CM6A-LR, and GFDL-
163 ESM4). For IPSL-CM6A-LR, temperatures are overestimated by the emulators throughout 1915-
164 2014. Close emulation of temperatures in abrupt-4xCO₂ does not guarantee close emulation for
165 the historical period (e.g. GFDL-ESM4, although using ERF from GFDL-CM4 likely introduces
166 some emulation error for GFDL-ESM4). Similarly, a relatively poor emulation of abrupt-4xCO₂
167 does not prohibit close emulation for the historical period (e.g. CNRM-CM6-1) (Figure S3).
168 These results are important because they show that there is no a priori way to know if an
169 AOGCM will be closely emulated.

170 The step-response emulator produces emulations with RMSEs broadly equivalent to or
171 less than emulations from EBM2. The treatment of time-varying feedbacks in the step-response
172 emulator (i.e., implicitly in the step-response function) differs from the treatment in EBM2 (i.e.,
173 based on ε) and may contribute to the good performance of the step-response emulator.

174 EBM3 performs better than EBM2 for abrupt-4xCO₂, which is expected given the
175 additional timescales resolved by the third layer which facilitates closer emulation of
176 temperatures during years 10-40 of the abrupt-4xCO₂ experiment, a period when the rate of
177 temperature increase weakens rapidly (Figure S3). However, the improvement of EBM3 over
178 EBM2 in the abrupt-4xCO₂ experiment does not consistently translate to the historical
179 experiment. Indeed there are three AOGCMs for which EBM2 has smaller RMSEs than EBM3
180 (HadGEM3-GC31-LL, MIROC6 and IPSL-CM6A-LR). EBM3, similar to EBM2, overestimates
181 temperatures for 2000-2014 in three of the eight AOGCMs and produces larger RMSEs than the
182 step-response emulator for some AOGCMs.



183
 184 **Figure 1.** Global mean temperature anomalies from a 1850-1900 baseline for CMIP6 AOGCMs.
 185 Changes in temperatures are forced by historical forcings during 1850-2014 and are shown for
 186 the period 1915-2014. RMSEs are calculated over 1915-2014.

187 3.2 Roles of different forcings for near-surface temperature change

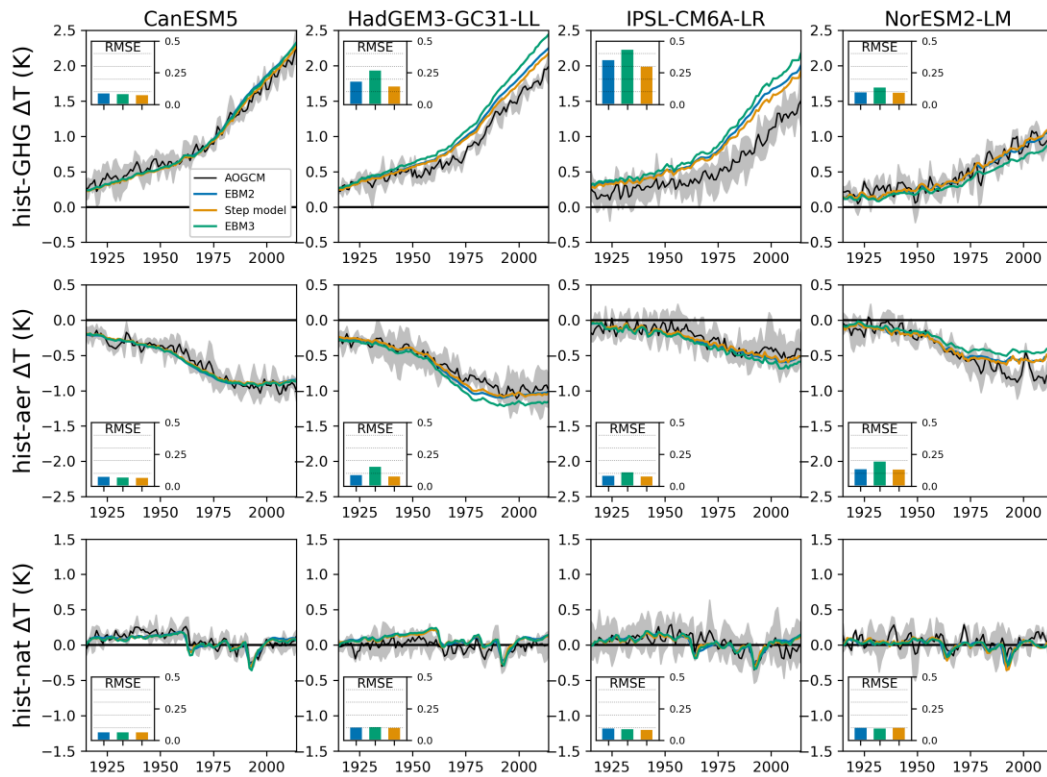
188 We used EBM2 to emulate the responses to historical greenhouse gas (hist-GHG),
 189 anthropogenic aerosol (hist-aer) and natural (hist-nat) forcings only. EBM2 was calibrated using
 190 abrupt-4xCO₂ simulations and the AOGCM projections are from DAMIP (Gillett et al., 2016)
 191 (Figure 2). We focus on two AOGCMs with relatively large errors in their emulations for the

192 historical period (HadGEM3-GC31-LL and IPSL-CM6A-LR), one AOGCM with relatively
193 small errors (CanESM5), and one AOGCM whose responses to GHG and aerosol forcings
194 contrast with the other AOGCMs (NorESM2-LM).

195 Although EBM2 was calibrated using abrupt-4xCO₂, errors predominantly arise from the
196 emulation of the response to GHG forcing; in part because GHGs have the largest ERF. The
197 EBM2 emulations overestimate the temperature increase due to GHGs for HadGEM3-GC31-LL
198 and IPSL-CM6A-LR.

199 Emulation of the temperature response to aerosol forcing is the largest source of error in
200 one climate model (NorESM2-LM). For HadGEM3-GC31-LL and IPSL-CM6A-LR, errors
201 associated with aerosol forcing offset errors associated with GHG forcing. This cancellation of
202 errors gives a spurious impression of better performance for the historical simulations. As shown
203 for the combined forcings (Figure 1), the step-response emulator produces closer emulations of
204 temperature for GHG forcing. For anthropogenic aerosol forcing, the step-response emulator
205 produces emulations of temperature very similar to EBM2.

206 Emulation of the temperature response to natural forcings is a small source of error and
207 the emulations are mostly within the spread of each AOGCM ensemble (Figures 2 and S4).
208 Although larger ensembles and longer simulations are required to robustly assess the emulated
209 response to volcanic forcing, thermal inertia of the EBM2 layers and allowance for rapid cloud
210 adjustments within RFMIP ERFs will likely have contributed to the close emulations for natural
211 forcings (Held et al., 2010; Gregory et al., 2016).



212

213 **Figure 2.** As Figure 1, except that temperature changes are forced by historical greenhouse gas
 214 (top row), anthropogenic aerosol (middle row), and natural (bottom row) forcings from RFMIP.

215 The AOGCM projections are from DAMIP.

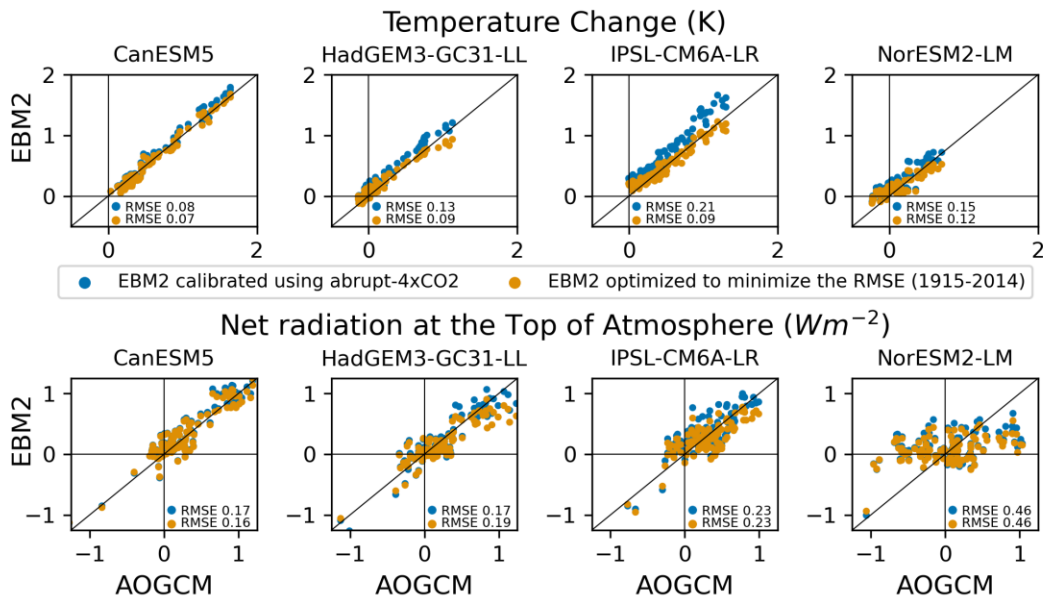
216 3.3 Alternative calibration of EBM2

217 To determine whether temperature emulations from EBM2 for the historical period can
218 be improved by changes to the fitted parameters alone, we apply optimization (Section 2.3) to
219 calibrate the λ and ε parameters (Figures 3 and S5).

220 This improves the emulations for all climate models. The greatest improvement occurs
221 during 1980-2014 and the emulation of temperature during this period is improved further if the
222 optimization is amended to minimize the RMSE specifically over this period. The spread in
223 emulated temperatures about the 1:1 line is mainly driven by the small AOGCM ensemble sizes
224 and is, therefore, similar for both EBM2 calibrations. Interannual variability is particularly large
225 for NorESM2-LM and the emulated temperatures have a low correlation with the AOGCM
226 temperatures for years prior to the 1980s when the climate response to forcing is relatively weak.

227 The emulations of the net radiation at the TOA (N) (Figure 3) show that close emulations
228 of near-surface temperature can be produced despite poor emulations of N . There is a large
229 spread in the emulations of N about the 1:1 line for all climate models. The emulation of N
230 during the late twentieth/early twenty-first century is poor for HadGEM3-GC31-LL and
231 emulated N has a weak correlation with its AOGCM for NorESM2-LM. Optimization does not
232 improve the emulation of N . There are small changes in emulated N for CanESM5 and
233 NorESM2-LM. The improved temperature emulations from the optimization method for
234 HadGEM3-GC31-LL come at the expense of poorer emulations of N . This result is important
235 because it demonstrates that climate model emulators can produce reasonable simulations of
236 near-surface temperature change, but the evolution of ocean heat uptake and TOA energy
237 imbalance is incorrect demonstrating limitations to physical interpretation.

238 We also used optimization to calibrate the λ and ε parameters separately for GHG and
239 aerosol forcing using the DAMIP experiments. The calibrated parameter values differ for the two
240 types of forcing (Table S3) and also vary when RMSE is minimized over different periods.



241
 242 **Figure 3.** Projected changes in global mean temperature (top row) and net radiation at the TOA
 243 (N) (bottom row). Each panel shows changes in the AOGCM (x-axis) against the EBM2
 244 emulation (y-axis). Each point represents an annual mean during 1915-2014.

245 3.4 Future near-surface temperature projections

246 We compare temperature emulations for the twenty-first century from EBM2 based on
 247 different methods for calibrating λ and ε (Figure 4). Results are shown for the AOGCMs where
 248 the most complete CMIP6 data are available. Results for other experiments are shown in Figure
 249 S6 and Table S1 describes the calibrations.

250 The performance of the abrupt-4xCO₂ calibration varies greatly between the AOGCMs
 251 (Figures 4a, b) and typically performs worse than the step-response emulator. The emulations of
 252 SSP2-4.5 deteriorate during the twenty-first century. The errors in the emulations are correlated
 253 with the magnitude of the forcing and peak near the end of the twenty-first century for total and
 254 GHG forcing and early in the twenty-first century for aerosol forcing.

255 Calibration by optimization of the λ and ε parameters over 1850-2100 (Figures 4c, d)
 256 yielded close emulations for all of the AOGCMs and across all experiments. Similarly close
 257 emulations were also achieved by minimizing the RMSE over 2015-2100 (not shown).
 258 Minimizing the RMSE for the later years of the projection, when the temperature anomalies are
 259 largest, is key.

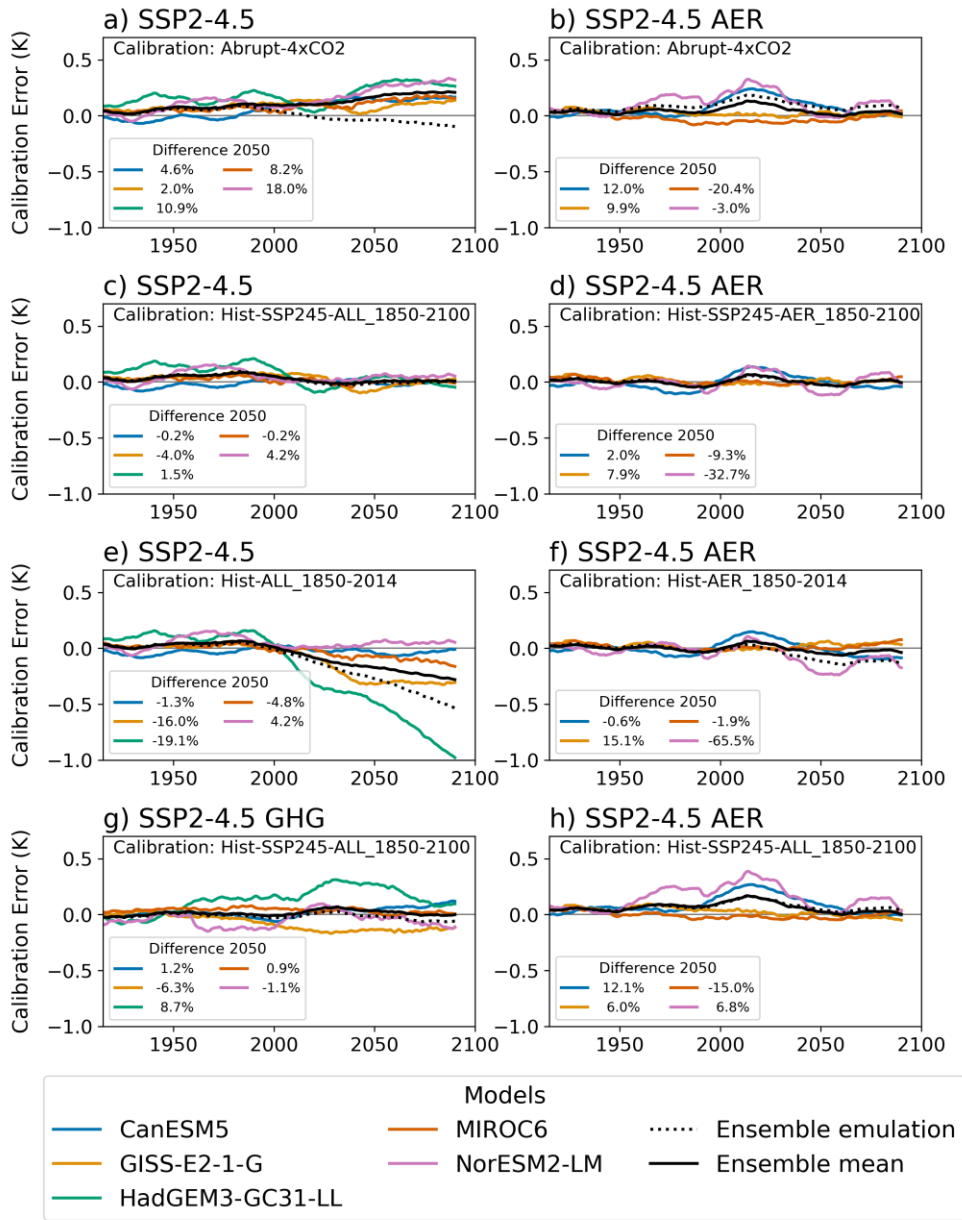
260 Emulations of temperature to 2100 based on optimizing the λ and ε parameters using the
261 1850-2014 period yields close emulations of temperature to 2014 but errors increase after the
262 calibration period (Figures 4e, f). Extending the calibration period from 1850-2014 to 1850-2040
263 (not shown) improves the emulation to 2040 but not always after 2040. Importantly, it does not
264 mitigate the risk of large emulation errors outside the calibration period and its impact varies
265 greatly between AOGCMs and between different experiments for the same AOGCM.

266 To investigate the impact of using a calibration from one experiment for a different
267 experiment, the “Hist-SSP245_1850-2100” calibration (which uses SSP2-4.5 all forcings) was
268 applied to the GHG only (SSP2-4.5-GHG) and the anthropogenic aerosol only (SSP2-4.5-AER)
269 experiments from DAMIP (Figures 4g, h). For both SSP2-4.5-GHG and SSP2-4.5-AER, the
270 error for the “Hist-SSP245_1850-2100” calibration is greater than for the Hist-SSP245-
271 GHG_1850-2100 and Hist-SSP245-AER_1850-2100 calibrations respectively. The impact also
272 varies between climate models and experiments in terms of the size of the impact and its
273 temporal behaviour. Similar results were also found for the high mitigation scenario SSP1-1.9
274 (Figure S7). Bespoke parameter calibrations for different ERF scenarios are necessary, therefore,
275 to achieve close emulations throughout 1850-2100. This result is important because it
276 demonstrates that emulator performance can be poor for out-of-sample predictions, yet there is
277 no clear a priori way to know if this will be the case. This poses a problem since some of the
278 value of emulators lies in their use for creating out-of-sample scenarios where AOGCM
279 simulations do not exist and cannot be readily performed.

280 The average of the emulations for individual climate models (Figure 4 “Ensemble
281 mean”) has relatively small RMSEs (except for the SSP2-4.5 1850-2014 calibration in Figure
282 4e). This is due, in part, to averaging of interannual variability across the ensemble of
283 emulations. Further, the ensemble mean generally has smaller RMSEs than an emulation in
284 which the ensemble mean ERF is used to emulate the ensemble temperature projection (Figure 4
285 “Ensemble emulation”).

286 Finally, while the optimization method yields unique parameter solutions there is a near
287 linear trade-off between the λ and ε parameters when minimizing the RMSE (demonstrated for
288 historical/SSP2-4.5 in Figure S2 and for the first 150 years of abrupt-4xCO₂ in Figure S8). In
289 EBM2, changes in the climate feedback parameter (λ) are compensated for by changes in the
290 efficacy of deep ocean heat uptake (ε) and the transient temperature response is largely

291 unchanged. This shows that optimized values for the λ and ε parameters may not be robust
 292 estimates of climate feedback or the AOGCM pattern effect and the physical interpretation of
 293 parameter value changes when optimizing the calibration is hindered by the linear trade-off
 294 between the λ and ε parameters.



295
 296 **Figure 4.** Differences between EBM2 emulations and AOGCM temperature projections. Rows
 297 show results for four calibrations of EBM2. Row B uses λ and ε parameter values which
 298 minimize the RMSE for temperatures during 1850-2100. Row C uses parameter values which
 299 minimize the RMSE during 1850-2014. Row D shows EBM2 calibrated to minimize the RMSE

300 during 1850-2100 for SSP2-4.5 and this calibration is used to emulate SSP2-4.5-GHG and SSP2-
301 4.5-AER. Annual means are smoothed using a 21-year moving average.

302 4 Discussion and conclusions

303 Our results show prediction errors in EBM2 for future global temperature projections
304 vary greatly between AOGCMs, forcings, time periods and methods of emulator calibration.

305 The EBM2 calibration using the abrupt-4xCO₂ experiment does not produce reliable
306 projections of historical warming for several AOGCMs. Although calibration of the λ and ε
307 parameters using optimization substantially reduces emulation errors for periods where an
308 AOGCM simulation is available, optimization of these parameters does not guarantee reliable
309 out-of-sample projections. Without an AOGCM projection for a given AOGCM and scenario, it
310 is not knowable if the EBM2 future projection will be reliable.

311 Close emulation of the historical period is not sufficient to guarantee reliable emulation
312 of future temperature changes (Figure 4; Nicholls et al., 2021). Opposing errors in the emulation
313 of GHG and aerosol forcings give a misleading impression of emulator performance. Many
314 climate model emulators do not reliably emulate AOGCM projections for high emissions
315 scenarios (Nicholls et al., 2021); our results suggest that strong mitigation scenarios may not be
316 reliably emulated.

317 How could the EBM2 emulator be changed to improve the out-of-sample emulations?
318 First, an efficacy factor could be introduced to account for asymmetry in AOGCM responses to
319 GHG and aerosol forcings. Second, EBM2 could be developed to incorporate variations in
320 climate feedbacks and the evolution of AOGCM pattern effects. Late twentieth-century warming
321 has been suppressed by changes in the observed sea surface temperature (SST) patterns and
322 associated cloud feedbacks (Andrews et al., 2018; Dong et al., 2021; Fueglistaler and Silvers,
323 2021). Future warming could be affected by changes in the pattern effect (Zhou et al., 2021).
324 Climate model simulations show that climate feedbacks weaken through time in response to
325 step-forcings and changes in feedbacks are associated with changes in SST patterns (e.g., Dong
326 et al., 2020; Dunne et al., 2020). Incorporating time-varying feedbacks in EBM2, however,
327 requires further research to distinguish forced changes in feedbacks from unforced climate noise
328 and to explicitly link global feedback changes to variations in SST patterns (e.g., using SST
329 anomalies for regions of tropical deep convection (Fueglistaler and Silvers (2021))).

330 EBM2 out-of-sample emulations could potentially be improved without changes to the
331 emulator. First, when available, larger AOGCM ensembles could be used to reduce the
332 contribution to emulation errors from chance. Second, more physically plausible parameter
333 tunings could be achieved by using optimization to jointly minimize RMSEs for temperature and
334 ocean heat flux (Dorheim et al., 2020). Our initial investigations minimizing RMSE for
335 temperature and N , however, showed that the emulation of historical temperatures was
336 substantially worse than when minimizing RMSE for temperature alone.

337 Emulations could also be improved through advances in the separation of forcing and
338 climate feedbacks in AOGCMs. We used the latest estimates of ERF derived from fixed-SST
339 simulations but substantial uncertainty in ERF remains (Forster et al., 2016; Dong et al., 2021).
340 Correcting for land warming in abrupt-4xCO₂ fixed-SST experiments increases the ERF
341 (Andrews et al., 2021) and leads to a weaker temperature response per unit forcing in EBM2. If
342 the abrupt-4xCO₂ ERF without corrections happens to be more underestimated than the
343 historical ERF, the historical EBM2 responses will be overestimated. Forcing estimates remain
344 dependent on the method used (Forster et al., 2013; Sherwood et al., 2015; Larson and Portmann,
345 2016; Fredriksen et al., 2021).

346 Our findings are relevant to observationally constrained climate model emulators aiming
347 to simulate real-world changes (e.g., Forster et al., 2021). Emulator structural errors and
348 uncertainties in inputs (e.g., ERF) are as relevant to real-world emulations as to emulations of
349 AOGCMs. Indeed, there are additional challenges. There is only one realization of past climate
350 and future climate is unknown. Observational large ensembles (McKinnon et al., 2017) could be
351 used to help characterize uncertainty in emulating past climate.

352 **Acknowledgments**

353 LSJ, ACM, TA and PMF were supported by the European Union’s Horizon 2020
354 programme under grant agreement No 820829 (CONSTRAIN). TA was supported by the Met
355 Office Hadley Centre Climate Programme funded by BEIS. CJS was supported by a joint
356 NERC-IIASA Collaborative Research Fellowship (NE/T009381/1). ACM was supported by The
357 Leverhulme Trust (PLP-2018-278). We acknowledge: the World Climate Research Programme
358 and its Working Group on Coupled Modeling for coordinating and promoting CMIP6; the
359 climate modeling groups for producing their model output; the Earth System Grid Federation
360 (ESGF) for archiving the data and providing access; and the funding agencies who support
361 CMIP6 and ESGF. We thank Nicholas Leach and an anonymous reviewer for their useful review
362 comments.

363 **Open Research**

364 The CMIP6 data were downloaded from the publicly available Earth System Grid
365 Federation archive at <https://esgf-node.llnl.gov/projects/cmip6/>. The R package for the three-
366 layer model (Cummins et al. 2020) was downloaded on 29 July 2021 from
367 <https://github.com/donaldcummins/EBM> and is available from
368 <https://doi.org/10.5281/zenodo.5217603>. Processed data produced for this paper are available on
369 Zenodo at <https://doi.org/10.5281/zenodo.6646804>.

370 **References**

- 371 Andrews, T., Gregory, J. M., & Webb, M. J. (2015). The dependence of radiative forcing and
372 feedback on evolving patterns of surface temperature change in climate models. *Journal of*
373 *Climate*, 28(4), 1630–1648. <https://doi.org/10.1175/JCLI-D-14-00545.1>.
- 374 Andrews, T., Gregory, J. M., Paynter, D., Silvers, L. G., Zhou, C., Mauritsen, T., Webb, M. J.,
375 Armour, K. C., Forster, P. M., & Titchner, H. (2018). Accounting for changing temperature
376 patterns increases historical estimates of climate sensitivity. *Geophysical Research Letters*, 45,
377 8490–8499. <https://doi.org/10.1029/2018GL078887>.
- 378 Andrews, T., Smith, C. J., Myhre, G., Forster, P. M., Chadwick, R., & Ackerley, D. (2021).
379 Effective radiative forcing in a GCM with fixed surface temperatures. *Journal of Geophysical*
380 *Research: Atmospheres*, 126, e2020JD033880. <https://doi.org/10.1029/2020JD033880>.
- 381 Armour, K. C., Bitz, C. M., & Roe, G. H. (2013). Time-Varying Climate Sensitivity from
382 Regional Feedbacks. *Journal of Climate*, 26(13), 4518–4534. [https://doi.org/10.1175/JCLI-D-12-](https://doi.org/10.1175/JCLI-D-12-00544.1)
383 00544.1.
- 384 Bloch-Johnson, J., Rugenstein, M., Stolpe, M. B., Rohrschneider, T., Zheng, Y., & Gregory, J.
385 M. (2021). Climate Sensitivity Increases Under Higher CO2 Levels Due to Feedback
386 Temperature Dependence. In *Geophysical Research Letters* (Vol. 48, Issue 4). Blackwell
387 Publishing Ltd. <https://doi.org/10.1029/2020GL089074>.
- 388 Colman, R., & Soldatenko, S. (2020). Understanding the links between climate feedbacks,
389 variability and change using a two-layer energy balance model. *Climate Dynamics*, 54, 3441–
390 3459, <https://doi.org/10.1007/s00382-020-05189-3>.
- 391 Cummins, D. P., Stephenson, D. B., and Stott, P. A. (2020). Optimal Estimation of Stochastic
392 Energy Balance Model Parameters. *Journal of Climate*, 33, 7909-7926, DOI: 10.1175/JCLI-D-
393 19-0589.1

- 394 Dong, Y., Armour, K. C., Zelinka, M. D., Proistosescu, C., Battisti, D. S., Zhou, C., & Andrews,
395 T. (2020). Intermodel spread in the pattern effect and its contribution to climate sensitivity in
396 CMIP5 and CMIP6 models. *Journal of Climate*, 33(18), 7755–7775.
397 <https://doi.org/10.1175/JCLI-D-19-1011.1>.
- 398 Dong, Y., Armour, K. C., Proistosescu, C., Andrews, T., Battisti, D. S., Forster, P. M., Paynter,
399 D., Smith, C. J., & Shiogama, H. (2021). Biased estimates of Equilibrium Climate Sensitivity
400 and Transient Climate Response derived from historical CMIP6 simulations. *Geophysical*
401 *Research Letters*. <https://doi.org/10.1029/2021GL095778>.
- 402 Dorheim, K., Link, R., Hartin, C., Kravitz, B., & Snyder, A. (2020). Calibrating Simple Climate
403 Models to Individual Earth System Models: Lessons Learned From Calibrating Hector. *Earth*
404 *and Space Science*, 7(11). <https://doi.org/10.1029/2019EA000980>.
- 405 Dunne, J. P., Winton, M., Bacmeister, J., Danabasoglu, G., Gettelman, A., Golaz, J. C., Hannay,
406 C., Schmidt, G. A., Krasting, J. P., Leung, L. R., Nazarenko, L., Sentman, L. T., Stouffer, R. J.,
407 & Wolfe, J. D. (2020). Comparison of Equilibrium Climate Sensitivity Estimates From Slab
408 Ocean, 150-Year, and Longer Simulations. *Geophysical Research Letters*, 47(16).
409 <https://doi.org/10.1029/2020GL088852>.
- 410 Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., & Taylor, K. E.
411 (2016). Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6) experimental
412 design and organization. *Geosci. Model Dev.*, 9, 1937–1958, doi:10.5194/gmd-9-1937-2016.
- 413 Forster, P. M., Andrews, T., Good, P., Gregory, J. M., Jackson, L. S., & Zelinka, M. (2013).
414 Evaluating adjusted forcing and model spread for historical and future scenarios in the CMIP5
415 generation of climate models. *J. Geophys. Res. Atmos.*, 118, 1139–1150.
416 <https://doi.org/10.1002/jgrd.50174>.
- 417 Forster, P. M., T. Richardson, A. C. Maycock, C. J. Smith, B. H. Samset, G. Myhre, T. Andrews,
418 R. Pincus, & M. Schulz (2016). Recommendations for diagnosing effective radiative forcing
419 from climate models for CMIP6, *J. Geophys. Res. Atmos.*, 121, 12,460–12,475,
420 doi:10.1002/2016JD025320.

- 421 Forster, P., T. Storelvmo, K. Armour, W. Collins, J. L. Dufresne, D. Frame, D. J. Lunt, T.
422 Mauritsen, M. D. Palmer, M. Watanabe, M. Wild, & H. Zhang (2021). The Earth's Energy
423 Budget, Climate Feedbacks, and Climate Sensitivity. In: *Climate Change 2021: The Physical*
424 *Science Basis. Contribution of Working Group I to the Sixth Assessment Report of the*
425 *Intergovernmental Panel on Climate Change* [Masson-Delmotte, V., P. Zhai, A. Pirani, S. L.
426 Connors, C. Péan, S. Berger, N. Caud, Y. Chen, L. Goldfarb, M. I. Gomis, M. Huang, K. Leitzell,
427 E. Lonnoy, J.B.R. Matthews, T. K. Maycock, T. Waterfield, O. Yelekçi, R. Yu and B. Zhou
428 (eds.)]. Cambridge University Press. In Press.
- 429 Fredriksen, H., Rugenstein, M., & Graversen, R. (2021). Estimating Radiative Forcing With a
430 Nonconstant Feedback Parameter and Linear Response. *Journal of Geophysical Research:*
431 *Atmospheres*, 126(24). <https://doi.org/10.1029/2020jd034145>.
- 432 Fueglistaler, S., & Silvers, L. G. (2021). The Peculiar Trajectory of Global Warming. *Journal of*
433 *Geophysical Research: Atmospheres*, 126(4). <https://doi.org/10.1029/2020JD033629>.
- 434 Geoffroy, O., Saint-Martin, D., Olivié, D. J. L., Voldoire, A., Bellon, G., & S. Tytéca, S.
435 (2013a). Transient Climate Response in a Two-Layer Energy-Balance Model. Part I: Analytical
436 Solution and Parameter Calibration Using CMIP5 AOGCM Experiments. *Journal of Climate*,
437 26, 1841-1857. doi: 10.1175/JCLI-D-12-00195.1.
- 438 Geoffroy, O., Saint-martin, D., Bellon, G., & Voldoire, A. (2013b). Transient Climate Response
439 in a Two-Layer Energy-Balance Model. Part II: Representation of the Efficacy of Deep-Ocean
440 Heat Uptake and Validation for CMIP5 AOGCMs. *Journal of Climate*, 26, 1859-1876. doi:
441 10.1175/JCLI-D-12-00196.1.
- 442 Gillett, N. P., Shiogama, H., Funke, B., Hegerl, G., Knutti, R., Matthes, K., Santer, B. D., Stone,
443 D., & Tebaldi, C. (2016). The Detection and Attribution Model Intercomparison Project
444 (DAMIP v1.0) contribution to CMIP6. *Geosci. Model Dev.*, 9, 3685–3697. doi:10.5194/gmd-9-
445 3685-2016.

- 446 Good, P., Gregory, J. M., & Lowe, J. A. (2011). A step-response simple climate model to
447 reconstruct and interpret AOGCM projections. *Geophysical Research Letters*, 38, L01703.
448 doi:10.1029/2010GL045208.
- 449 Good, P., Lowe, J. A., Andrews, T., Wiltshire, A., Chadwick, R., Ridley, J. K., Menary, M. B.,
450 Bouttes, N., Dufresne, J. L., Gregory, J. M., Schaller, N., & Shiogama, H. (2015). Nonlinear
451 regional warming with increasing CO₂ concentrations. *Nature Climate Change*, 5(2), 138–142.
452 doi.org/10.1038/nclimate2498.
- 453 Gregory, J. M., Andrews, T., & Good, P. (2015). The inconstancy of the transient climate
454 response parameter under increasing CO₂. *Phil. Trans. R. Soc. A* 373: 20140417.
455 <http://dx.doi.org/10.1098/rsta.2014.0417>.
- 456 Gregory, J. M., Andrews, T., Good, P., Mauritsen, T., & Forster, P. M. (2016). Small
457 global-mean cooling due to volcanic radiative forcing. *Clim. Dyn.*, 47, 3979–3991. DOI
458 10.1007/s00382-016-3055-1.
- 459
- 460 Held, I. M., Winton, M., Takahashi, K., Delworth, T., Zeng, F., & Vallis, G. K. (2010). Probing
461 the Fast and Slow Components of Global Warming by Returning Abruptly to Preindustrial
462 Forcing. *Journal of Climate*, 23, 2418-2427. Doi: 10.1175/2009JCLI3466.1.
- 463 Larson, E. J. L., & Portmann, R. W. (2016). A Temporal Kernel Method to Compute Effective
464 Radiative Forcing in CMIP5 Transient Simulations. *Journal of Climate*, 29, 1497–1509.
465 <https://doi.org/10.1175/JCLI-D-15-0577.1>.

- 466 Lee, J. Y., J. Marotzke, G. Bala, L. Cao, S. Corti, J. P. Dunne, F. Engelbrecht, E. Fischer, J. C.
467 Fyfe, C. Jones, A. Maycock, J. Mutemi, O. Ndiaye, S. Panickal, & T. Zhou (2021). Future
468 Global Climate: Scenario-Based Projections and Near-Term Information. In: *Climate Change*
469 *2021: The Physical Science Basis. Contribution of Working Group I to the Sixth Assessment*
470 *Report of the Intergovernmental Panel on Climate Change* [Masson-Delmotte, V., P. Zhai, A.
471 Pirani, S. L. Connors, C. Péan, S. Berger, N. Caud, Y. Chen, L. Goldfarb, M. I. Gomis, M.
472 Huang, K. Leitzell, E. Lonnoy, J. B. R. Matthews, T. K. Maycock, T. Waterfield, O. Yelekçi, R.
473 Yu and B. Zhou (eds.)]. Cambridge University Press. In Press.
- 474 McKinnon, K. A., Poppick, A., Dunn-Sigouin, E., & Deser, C. (2017). An “Observational Large
475 Ensemble” to Compare Observed and Modeled Temperature Trend Uncertainty due to Internal
476 Variability. *Journal of Climate*, 30, 7585–7598. <https://doi.org/10.1175/JCLI-D-16-0905.1>.
- 477 Modak, A., & Mauritsen, T. (2021). The 2000–2012 global warming hiatus more likely with a
478 low climate sensitivity. *Geophysical Research Letters*, 48, e2020GL091779.
479 <https://doi.org/10.1029/2020GL091779>.
- 480 Nicholls, Z. R. J., Meinshausen, M., Lewis, J., Gieseke, R., Dommenges, D., Dorheim, K., Fan, S.,
481 Fuglestad, J. S., Gasser, T., Golube, U., Goodwin, P., Hartin, C., P. Hope, A., Kriegler,
482 E., J. Leach, N., Marchegiani, D., A. McBride, L., Quilcaille, Y., Rogelj, J., & Xie, Z. (2020).
483 Reduced Complexity Model Intercomparison Project Phase 1: Introduction and evaluation of
484 global-mean temperature response. *Geoscientific Model Development*, 13(11), 5175–5190.
485 <https://doi.org/10.5194/gmd-13-5175-2020>.
- 486 Nicholls, Z., Meinshausen, M., Lewis, J., Corradi, M. R., Dorheim, K., Gasser, T., Gieseke, R.,
487 Hope, A. P., Leach, N. J., McBride, L. A., Quilcaille, Y., Rogelj, J., Salawitch, R. J., Samset, B.
488 H., Sandstad, M., Shiklomanov, A., Skeie, R. B., Smith, C. J., Smith, S. J., Su, X., Tsutsui, J.,
489 Vega-Westhoff, B., & Woodard, D. L. (2021). Reduced complexity Model Intercomparison
490 Project Phase 2: Synthesizing Earth system knowledge for probabilistic climate projections.
491 *Earth's Future*, 9, e2020EF001900. <https://doi.org/10.1029/2020EF001900>.

- 492 Pincus, R., Forster, P. M., & Stevens, B. (2016), The Radiative Forcing Model Intercomparison
493 Project (RFMIP): experimental protocol for CMIP6. *Geosci. Model Dev.*, 9, 3447–3460.
494 doi:10.5194/gmd-9-3447-2016.
- 495 Riahi, K., Schaeffer, R., Arango, J., Calvin, K., Guivarch, C., Hasegawa, T., et al. (2022):
496 Mitigation pathways compatible with long-term goals. In IPCC, 2022: *Climate Change 2022:
497 Mitigation of Climate Change. Contribution of Working Group III to the Sixth Assessment
498 Report of the Intergovernmental Panel on Climate Change* [P.R. Shukla, J. Skea, R. Slade, A. Al
499 Khourdajie, R. van Diemen, D. McCollum, M. Pathak, S. Some, P. Vyas, R. Fradera, M.
500 Belkacemi, A. Hasija, G. Lisboa, S. Luz, J. Malley, (eds.)]. Cambridge University Press,
501 Cambridge, UK and New York, NY, USA. doi: 10.1017/9781009157926.005
- 502 Rohrschneider, T., Stevens, B., & Mauritsen, T. (2019). On simple representations of the climate
503 response to external radiative forcing. *Climate Dynamics*, 53(5–6), 3131–3145.
504 <https://doi.org/10.1007/s00382-019-04686-4>.
- 505 Rugenstein, M. A. A., Caldeira, K., & Knutti, R. (2016). Dependence of global radiative
506 feedbacks on evolving patterns of surface heat fluxes. *Geophysical Research Letters*, 43(18),
507 9877–9885. <https://doi.org/10.1002/2016GL070907>.
- 508 Rugenstein, M., Bloch-Johnson, J., Gregory, J., Andrews, T., Mauritsen, T., Li, C., Frölicher, T.
509 L., Paynter, D., Danabasoglu, G., Yang, S., Dufresne, J. L., Cao, L., Schmidt, G. A., Abe-Ouchi,
510 A., Geoffroy, O., & Knutti, R. (2020). Equilibrium Climate Sensitivity Estimated by
511 Equilibrating Climate Models. *Geophysical Research Letters*,
512 47(4).<https://doi.org/10.1029/2019GL083898>.
- 513 Senior, C. A., & Mitchell, J. F. B. (2000). The time-dependence of climate sensitivity.
514 *Geophysical Research Letters*, 27(17), 2685–2688. <https://doi.org/10.1029/2000GL011373>.
- 515 Sherwood, S. C., Bony, S., Boucher, O., Bretherton, C., Forster, P. M., Gregory, J. M., &
516 Stevens, B. (2015). Adjustments in the forcing-feedback framework for understanding climate
517 change. *Bulletin of the American Meteorological Society*, 96(2), 217–228.
518 <https://doi.org/10.1175/BAMS-D-13-00167.1>.

- 519 Smith, C. J., Kramer, R. J., Myhre, G., Alterskjær, K., Collins, W., Sima, A., et al. (2020).
520 Effective radiative forcing and adjustments in CMIP6 models. *Atmos. Chem. Phys.*, 20, 9591–
521 9618, <https://doi.org/10.5194/acp-20-9591-2020>.
- 522 Smith, C. J., Harris, G. R., Palmer, M. D., Bellouin, N., Collins, W., Myhre, G., Schulz, M.,
523 Golaz, J.-C., Ringer, M., Storelvmo, T., & Forster, P. M. (2021). Energy Budget Constraints on
524 the Time History of Aerosol Forcing and Climate Sensitivity. *Journal of Geophysical Research:*
525 *Atmospheres*, 126, e2020JD033622. <https://doi.org/10.1029/2020JD033622>.
- 526 Stevens, B., Sherwood, S. C., Bony, S., & Webb, M. J. (2016). Prospects for narrowing bounds
527 on Earth's equilibrium climate sensitivity, *Earth's Future*, 4, 512–522.
528 [doi:10.1002/2016EF000376](https://doi.org/10.1002/2016EF000376).
- 529 Winton, M., Takahashi, K., & Held, I. M. (2010). Importance of Ocean Heat Uptake Efficacy to
530 Transient Climate Change. *Journal of Climate*, 23, 2333–2344, DOI: 10.1175/2009JCLI3139.1.
- 531 Zhou, C., Zelinka, M. D., Dessler, A. E., & Wang, M. (2021). Greater committed warming after
532 accounting for the pattern effect. *Nature Climate Change*, 11(2), 132–136.
533 <https://doi.org/10.1038/s41558-020-00955-x>.