*Note that all line numbers in our responses refer to the version with tracked comments included in this document. Line numbers from the reviewers refer to the original manuscript.*

Reviewer #1

We thank the reviewer for their comments. Below, we detail our responses.

*64-65. As a single number to quantify the spread, the standard deviation would also be helpful.*

We have added 5-95% confidence intervals throughout the paper.

*66. Why do you use only a single decade, rather than all the data, for instance by dividing the dataset into two or using regression (cf Barnes and Barnes, 2015, 10.1175/JCLI-D-15-0032.1)? A single decade would be less precise. You could estimate the statistical uncertainty incurred from the control run.*

We calculate ECS using this approach because this is the way most ECS calculations based on the 20[th]-century observational record are done. Thus, our results can therefore directly provide insight into the impact of variability in the observational estimates of ECS.

The reviewer is correct that using more than a decade might affect the results. If one used the difference between the averages of the first and last 20 years, the range in $\lambda$ declines from 0.87 W/m$^2$/K to 0.48 W/m$^2$/K. Using longer averaging periods does not further decrease the range. We now mention this in the paper (line 67).

*118. It would be useful to remark here that 16 years is chosen to match the CERES dataset, because that was mentioned some lines above (103-104), where it appears actually to be 17 years and 5 months long.*

We have added a statement that the segmentation of the data is done to match the CERES record (line 134). We have also updated the paper to segment the data into 17-year segments to more closely match CERES.

*119, 196. Why are monthly anomalies used here, rather than annual? Does it make a difference?*

We do this to facilitate the comparison with the CERES regressions, which also uses monthly data. The reason most analyses with CERES data are done with monthly data is because using annual data means there's only 17 data points, and the uncertainties end up being very large. Issues involved in annual vs. monthly regressions are discussed in some detail in Forster (2016, 10.1146/annurev-earth-060614-105156).

*167. Again, the standard deviation would be helpful, and could be compared with lines 64-65.*

Added.

*173, 175. You could give standard errors of the mean for each of these two numbers, and judge the significance of their difference.*

We have added the 5-95% confidence intervals to all of these numbers.

*174, 175. "analysis" and "calculated" - by what method? From the slope of R against Delta T?*

We have clarified the text that we use the method of Gregory et al. (2004), where annual average R is regressed against T, and the slope of the curve is an estimate of λ or Θ (line 194)

*204. "agrees" in what sense?*

We have changed the sentence to read: "We find that the 15 models whose average short-term Θ falls within the uncertainty of Θ estimated from CERES observations have ECS values ranging from 2.0-3.9 K, with an average of 2.9 K." (line 247)

*218. I would say that this is "one source" of the spread, which is not eliminated, but only reduced, by using Theta instead.*

We believe that this sentence is phrased correctly. The spread in our estimate from the ensemble is due to the construction of the energy-balance equation. Unlike observational analyses, we know everything else perfectly. Using our revised energy balance equation does not completely solve the problem, but it is an improvement.

*233. Why is this material an appendix, rather than being incorporated in the main text?*

We felt that this material would not be interesting to most readers, so we put it in the appendix. In retrospect, perhaps that was a bad decision. At this stage in the paper's review cycle, we hesitate to move material around. We can, however, if the reviewer or editor insists.

Reviewer #2

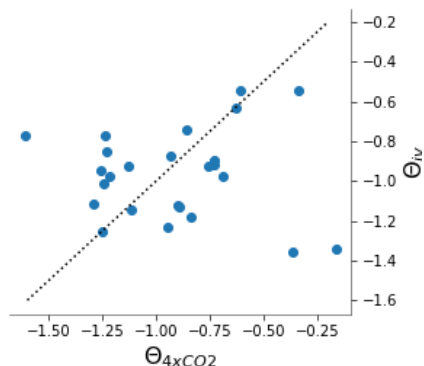We thank the reviewer for their comments. In this document, we detail our responses.

*1) It would be helpful to provide a little more physical motivation for the choice of tropical 500 hPa temperature. I see some good reasons why mid-tropospheric temperature should work better (e.g., it should scale better with LR, WV and LW cloud feedbacks), but I don't think this was discussed anywhere. Why use tropical temperature rather than global-mean? Is there a physical rationale, or did this simply work better in MPI-ESM?*

*Also, although mid-tropospheric temperature clearly works better for the overall feedback, I expect the scaling with Ta might actually be a worse choice for some individual feedback processes (e.g. surface albedo, marine low cloud). This might be worth discussing briefly.*

A: To address this, we have added a paragraph to the paper beginning on line 221.
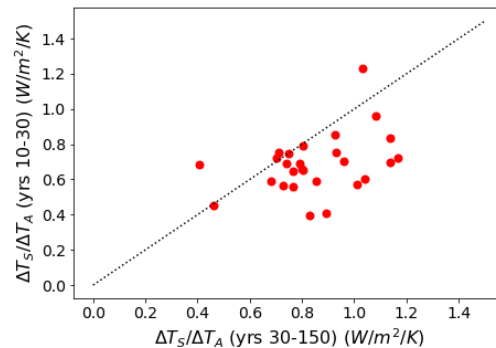
*2) A key result is that the revised feedback parameter theta more accurately estimates the "true" feedback strength under CO2 forcing. This is shown to be the case in MPI-ESM (L172-176). However, does this hold for CMIP5 models in general? I.e., do the values of theta estimated in control runs correlate well with those in 4xCO2?*

A: This is not a claim we make in the paper, although one might infer it from the MPI model. Indeed, there is *some* correlation between short-term and long-term theta in the CMIP5 ensemble, as seen here:



Caption: Scatter plot of $\Theta_{4xCO2}$ vs. $\Theta_{control}$ from the CMIP5 ensemble. Each point represents values from model.

However, because of the outlier models, the relation is hard to interpret and we have not pursued this "emergent constraint" approach in our estimate of ECS using our revised framework [Dessler and Forster (2018, February 6). An estimate of equilibrium climate sensitivity from interannual variability. Retrieved from eartharxiv.org/4et67].

We have added a short statement to the paper to reflect this on line 260: "It may also be possible to use the relation between short-term and long-term $\Theta$ as an emergent constraint to convert short-term observations to the long-term response. There is some scatter in the relation in the CMIP5 ensemble, however, so more analysis of how these relate is likely required before ECS can be constrained in this way."

*Relatedly, I would also suggest adding the correlation between R and Ta in CMIP5 piControl to Fig. 4, as additional bars in a different color.*

A: We have done that.

*3) One important issue that isn't discussed in the paper is that the "pattern effect" doesn't simply go away with the improved relationship; rather, it shifts from the feedback parameter to the Ts/Ta term. This isn't a problem, but the way the paper is currently written, some readers might get that impression.*

A: We have added a sentence discussing this: "Thus, the pattern effect's impact on ECS calculations shifts from $\lambda$ in the traditional framework to the $\Delta T_S / \Delta T_A$ term in Eq. 4." (line 217)

*So if most of the curvature in the relationship between radiative response and temperature goes away with the revised framework (Fig. 6), I expect there must be some curvature in the Ta versus Ts relationship in 4xCO2 runs. Can the authors confirm this?*

Confirmed.



Caption. Scatterplot of slope of $\Delta T_S$ vs. $\Delta T_A$ in CMIP5 abrupt4xCO2 runs. Each point represents one model. The dotted line is the 1:1 line. The subscripts (10-30, 30-150) indicate the years of the run from which the slopes are calculated.

We've added a sentence to the paper mentioning that there is curvature in $T_A$ vs $T_S$ relation: "The lack of curvature in the $\Theta$ calculations means there is curvature in the relation between $T_A$ and $T_S$ in the models." (line 216)

*4) I expect the Ts/Ta ratio cannot be reliably estimated from historical runs in the presence of large variability (for the same reason that lambda cannot be reliably estimated - because of the pattern effect). So we must rely on models to estimate this ratio under future global warming, meaning that it will be important to understand how future patterns of surface warming will develop. I suggest the authors discuss this briefly, for example in the conclusions.*

We have added a sentence to the paper mentioning this point: "This also emphasizes the need to improve our understanding of the factors that control $\Delta T_S/\Delta T_A$, as well as how future patterns of surface warming will evolve." (line 218)

*Other minor comments:*

*I suggest using colors in Fig. 6, rather than dark grey and black.*

Done

*L223: Cite Andrews and Webb 2018 - For future reference, it would be useful to mention the value of theta estimated from observations (horizontal dashed bar in Fig. 7a).*

Done.

**The influence of internal variability on Earth's energy balance framework and implications for estimating climate sensitivity**

Andrew E. Dessler[1]*, Thorsten Mauritsen[2], Bjorn Stevens[2]

[1] Dept. of Atmospheric Sciences, Texas A&M University, College Station, TX 77843
[2] Max Planck Institute for Meteorology, Bundesstraße 53, 20146 Hamburg, Germany

*Correspondence to: adessler@tamu.edu, 979-862-1427

1

**Abstract:** Our climate is constrained by the balance between solar energy absorbed by the
Earth and terrestrial energy radiated to space.  This energy balance has been widely used to
infer equilibrium climate sensitivity (ECS) from observations of 20[th]-century warming.  Such
estimates yield lower values than other methods and these have been influential in pushing
down the consensus ECS range in recent assessments.  Here we test the method using a 100-
member ensemble of the MPI-ESM1.1 climate model simulations of the period 1850-2005 with
known forcing.  We calculate ECS in each ensemble member using energy balance, yielding
values ranging from 2.1 to 3.9 K.  The spread in the ensemble is related to the central
hypothesis in the energy budget framework: that global average surface temperature
anomalies are indicative of anomalies in outgoing energy (either of terrestrial origin or reflected
solar energy).  We find that assumption is not well supported over the historical temperature
record in the model ensemble or more recent satellite observations.  We find that framing
energy balance in terms of 500-hPa tropical temperature better describes the planet's energy
balance.

23 **The problem**

24 When an energy imbalance is imposed, such as by adding a greenhouse gas to the atmosphere,

25 the climate will shift in such a way to eliminate the energy imbalance.   This process is

26 embodied in the traditional linearized energy balance equation:

27 $\qquad R = F + \lambda\,T_s$ (1)

28 where the forcing F is an imposed energy imbalance, $T_S$ is the global average surface

29 temperature, and $\lambda$ relates changes in $T_S$ to a change in net top-of-atmosphere (TOA) flux

30 (Gregory et al., 2002; Dessler and Zelinka, 2014).  R is the resulting TOA flux imbalance from the

31 combined forcing and response.  All quantities are deviations from an equilibrium base state,

32 usually the pre-industrial climate. Equilibrium climate sensitivity (hereafter ECS, the equilibrium

33 warming in response to a doubling of $CO_2$) is equal to $-F_{2xCO2}/\lambda$, where $F_{2xCO2}$ is the forcing from

34 doubled $CO_2$.

35 Many investigators (e.g., Gregory et al., 2002; Annan and Hargreaves, 2006; Otto et al., 2013;

36 Lewis and Curry, 2015; Aldrin et al., 2012; Skeie et al., 2014; Forster, 2016) have used Eq. 1

37 combined with estimates of R, F, and $T_s$ to estimate $\lambda$:

38 $\qquad \lambda = \Delta(R-F)/\Delta T_s$ (2)

39 where $\Delta$ indicates the change between the start of the historical period (usually the mid to late

40 nineteenth century) and a recent period.  These calculations result in values of $\lambda$ near

41 -2 $W/m^2/K$ and appear to rule out ECS larger than ~4 K (Stevens et al., 2016).  The substantial

42 likelihood of an ECS below 2 K implied by these calculations led the IPCC Fifth Assessment

43 Report to extend their lower bound on *likely* values of ECS to 1.5 K (Collins et al., 2013).

44 We test this energy balance methodology through a perfect model experiment consisting of an

45 analysis of a 100-member ensemble of runs of the MPI Earth System Model, MPI-ESM1.1.  This

46 is the latest coupled climate model from the Max Planck Institute for Meteorology and consists

47 of the ECHAM6.3 atmosphere and land model coupled to the MPI-OM ocean model. The

48 atmospheric resolution is T63 spectral truncation, corresponding to about 200 km, with 47

49  vertical levels, whereas the ocean has a nominal resolution of about 1.5 degrees and 40 vertical

50  levels. MPI-ESM1.1 is a bug-fixed and improved version of the MPI-ESM used during CMIP5

51  (Giorgetta et al., 2013) and nearly identical to the MPI-ESM1.2 (Mauritsen et al., 2018) model

52  being used to provide output to CMIP6, except that the historical forcings are from the MPI-

53  ESM.

54  Each of the 100 members simulates the years 1850-2005 (Fig. 1) and use the same evolution of

55  historical natural and anthropogenic forcings.  The members differ only in their initial

56  conditions —each starts from a different state sampled from a 2000-year control simulation.

57  We calculate effective radiative forcing F for the ensemble by subtracting top-of-atmosphere

58  flux R in a run with climatological sea surface temperatures (SSTs) and a constant pre-industrial

59  atmosphere from average R from an ensemble of three runs using the same SSTs but the time-

60  varying atmospheric composition used in the historical runs (Hansen et al., 2005; Forster et al.,

61  2016).  The three-member ensemble begins with perturbed atmospheric states.  We estimate

62  $F_{2xCO2}$ using the same approach in a set of fixed SST runs in which $CO_2$ increases at 1% per year,

63  which yields a $F_{2xCO2}$ value of 3.9 W/m$^2$.

64  We calculate $\lambda$ using Eq. 2 for each ensemble member, producing values ranging from -1.88 to

65  -1.01 W/m$^2$/K (5-95% range -1.63 to -1.17 W/m$^2$/K), with an ensemble median of -1.43 W/m$^2$/K

66  (Fig. 2a).  In this calculation, $\Delta$(R-F) and $\Delta T_S$ are the average difference between the first and last

67  decade of each run.  The spread in $\lambda$ depends to some extent on how the calculation is set up

68  — if one used the difference between the averages of the first and last 20 years, for example,

69  the range in $\lambda$ declines from 0.87 W/m$^2$/K to 0.48 W/m$^2$/K.  Using longer averaging periods does

70  not further decrease the range.

71  We also calculate ECS = $-F_{2xCO2}/\lambda$ for each ensemble member, producing values ranging from

72  2.08 to 3.87 K (5-95% range 2.39 to 3.34 K) (Fig. 2b), with an ensemble median of 2.72 K.  Thus,

73  our analysis shows that $\lambda$ and ECS estimated from the historical record can vary widely simply

74  due to internal variability. Given that we have only a single realization of the 20[th] century, we

75  should not consider estimates based on the historical period to be precise — even with perfect

76  observations. This supports previous work that also emphasized the impact of internal

**Deleted:** average

**Deleted:** average

**Deleted:** 6

**Deleted:** ¶
With respect to precision of the estimates

82  variability on estimates of $\lambda$ and ECS (Huber et al., 2014; Andrews et al., 2015; Zhou et al., 2016;

83  Gregory and Andrews, 2016).

84  Previous researchers have questioned whether the historical record provides an accurate

85  measure of $\lambda$ and ECS, and we can check this by comparing the ensemble values to ECS

86  estimates from a $2xCO_2$ run of the MPI-ESM1.2, which is physically very close to MPI-ESM1.1.

87  An abrupt $2xCO_2$ run yields an ECS of 2.93 K in response to an abrupt doubling of $CO_2$

88  (estimated by regressing years 100-1000 of a 1000-year run) — 8% larger than the ensemble

89  median. This is in line with the 10% difference in ECS estimated by Mauritsen and Pincus (2017)

90  to arise from the average CMIP5 model time-dependent feedback, but smaller than suggested

91  in other recent studies of ECS in transient climate runs (e.g., Armour, 2017; Proistosescu and

92  Huybers, 2017).

93  Thus, there are a number of issues that need to be considered when interpreting estimates of $\lambda$

94  and ECS derived from the historical period.  In addition to the precision and accuracy issues

95  discussed above, it also includes the large and evolving uncertainty in forcing over the 20[th]

96  century (Forster, 2016), different forcing efficacies of greenhouse gases and aerosols (Shindell,

97  2014; Kummer and Dessler, 2014), and geographically incomplete or inhomogeneous

98  observations (Richardson et al., 2016).

99  **Why are estimates using the traditional energy balance approach imprecise?**

100 In this section, we explain the physical process by which internal variability leads to the large

101 spread in $\lambda$ and ECS estimated from the ensemble.  We begin by observing that Eqs. 1 and 2

102 parameterize R-F in terms of global average surface temperature, $T_S$.  In model runs with strong

103 forcing driving large warming, such as abrupt $4xCO_2$ simulations, there is indeed a strong

104 correlation between these variables (e.g., Gregory et al., 2004).  However, because R-F in such

105 runs is dominated by a monotonic trend, correlations will exist with any geophysical field that

106 also exhibits a monotonic trend, regardless of whether there is a physical connection between

107 the fields. Thus, one should not take the correlation between R-F and $T_S$ in these runs as

108 proving causality.

---

Deleted: average

Deleted:  and the changes between the MPI-ESM1.1 and MPI-ESM1.2 are not believed to be important for its climate sensitivity

Deleted: 6

Deleted: average

115  If $T_S$ is a good proxy for the response R-F, we would expect to also see a correlation in

116  measurements dominated by interannual variations. Observational data allow us to test this

117  hypothesis.  We use observations of R from the Clouds and the Earth's Radiant Energy System

118  (CERES) Energy Balanced and Filled product (ed. 4) (Loeb et al., 2009), which cover the period

119  March 2000 to July. 2017. Our sign convention throughout the paper is that downward fluxes

120  are positive.  Temperatures come from the European Centre for Medium Range Weather

121  Forecasts (ECMWF) Interim Re-Analysis (ERAi) (Dee et al., 2011).  We assume forcing changes

122  linearly over this time period and account for it by detrending $\Delta$R and $\Delta$T anomaly time series

123  using a linear least-squares fit to remove the long-term trend.

124  These data show that $\Delta$R is poorly correlated with $\Delta T_s$ in response to interannual variability (Fig.

125  3a), as has been noted many times in the literature; see, e.g., Sect. 5 of Forster (2016).  In

126  particular, the low correlation coefficient tells us that $\Delta T_S$ explains little of the variance in $\Delta$R.

127  Using explicit estimates of forcing or other temperature datasets (e.g., MERRA-2) yield the

128  same result.

129  GCMs that submitted output to the 5[th] phase of the Coupled Model Intercomparison Project

130  (CMIP5) (Taylor et al., 2012) also show this poor correlation.  To demonstrate this, we have

131  calculated the correlation coefficient between $\Delta T_S$ and $\Delta$R in CMIP5 pre-industrial control runs

132  (these are runs for which forcing F = 0).  To facilitate comparison with the CERES data, as well as

133  avoid any issues with long-term drift in the control runs, we break each run into 17-year                    Deleted: 6

134  segments to match the length of the CERES data and calculate the correlation coefficient of

135  monthly anomalies of $\Delta$R and $\Delta T_S$ for each segment. Fig. 4 shows that the correlation between

136  $\Delta$R and $\Delta T_S$ in the models is similar to that from the CERES analysis.

137  Recent work provides an explanation: the response of $\Delta$(R-F) to a particular $\Delta T_S$ is determined

138  not only by the global average magnitude, but also by the pattern of warming (Armour et al.,

139  2013; Andrews et al., 2015; Gregory and Andrews, 2016; Zhou et al., 2016, 2017; Andrews and

140  Webb, 2018). During El Nino cycles that dominate the observations in Fig. 3, the spatial pattern

141  of warm and cool regions changes, leading to responses in $\Delta$(R-F) that do not scale cleanly with

142  $\Delta T_S$ — something Stevens et al. (2016) refer to as "pattern effects"                    Deleted: (2016)

6

145    To demonstrate how this also generates the spread in $\lambda$ in the model ensemble (Fig. 2a), we

146    calculate the local response $\lambda_r$ in three equal-area regions (90°S-19.4°S, 19.4°S-19.4°N, 19.4°N-

147    90°N).  We define $\lambda_r$ as the regional analog to $\lambda$ (Eq. 2):

148         $\lambda_r = \Delta(R-F)_r/\Delta T_{S,r}$                                    (3)

149    where the "r" subscript indicates a regional average value.

150    We find that $\lambda_r$ varies between the regions (Fig. 5). This means that different ensemble

151    members with similar global average $\Delta T_S$ but different patterns of surface warming produce

152    different values of global average $\Delta(R-F)$, thereby leading to spread in the estimated $\lambda$ among

153    the ensemble members.  We also see strong variability in $\lambda_r$ within each region, suggesting that

154    how the warming is distributed within the region also drives some of the spread in estimated $\lambda$

155    in the ensemble.

156    This explanation is consistent with analyses showing that $\lambda$ changes during transient runs as the

157    pattern of surface temperature evolves (Senior and Mitchell, 2000; Armour et al., 2013;

158    Andrews et al., 2015; Gregory and Andrews, 2016; Stevens et al., 2016).  In our model

159    ensemble, however, the pattern changes are caused by internal variability rather than differing

160    regional heat capacities that cause some regions to warm more slowly than others during

161    forced warming.

162    **A better way to describe energy balance**

163    Our analysis demonstrates limitations of the conventional energy balance framework (Eq. 1). It

164    has been previously noted that $\Delta R$ correlates better with tropospheric temperatures than $\Delta T_S$

165    (Murphy, 2010; Spencer and Braswell, 2010; Trenberth et al., 2015). Recent analyses have also

166    stressed the importance of atmospheric temperatures — through its influence on lapse rate —

167    as providing a fundamental control on the planet's energy budget (Zhou et al., 2016; Ceppi and

168    Gregory, 2017).  Based on this, we test a new energy balance framework constructed using the

169    temperature of the tropical atmosphere:

170    R - F = Θ T_A            (4)

171  where $T_A$ is the tropical average (30°N-30°S) 500-hPa temperature and Θ relates this quantity to

172  R-F.  R and F are the same global average quantities they were in equation 1.  ECS can be

173  expressed in terms of Θ:

174    $$ECS = -\frac{\Delta F_{2\times CO2}}{\Theta}\frac{\Delta T_S}{\Delta T_A}$$      (5)

175  where $\Delta T_S$ and $\Delta T_A$ are the equilibrium changes in these quantities in response to doubled $CO_2$.

176  The CMIP5 ensemble average ratio $\Delta T_S/\Delta T_A$ is 0.86±0.10 (±1σ), where Δ represents the average

177  difference between the first and last decades of the abrupt $4xCO_2$ runs.

178  Support for Eq. 4 can be found in the observations: ΔR shows a tighter correlation with $\Delta T_A$ than

179  with $\Delta T_S$ in observations (Figs. 3a vs. 3b).  CMIP5 models also show this (Fig. 4).  Given that the

180  slope of these plots can be taken as estimates of Θ and λ, the tighter correlation leads to more

181  accurate estimates of Θ than λ, both in absolute and relative terms.

182  Turning to the model ensemble, we next demonstrate that Θ is a more precise metric than λ.

183  We do this by calculating Θ [= Δ(R-F)/$\Delta T_A$] in each ensemble member, yielding values ranging

184  from -1.18 to -0.89 W/m²/K (5-95% range -1.16 to -0.92 W/m²/K), with an ensemble median of

185  -1.04 W/m²/K (Fig. 2a).  There is clearly less variability in Θ among the ensemble members than

186  for λ.  This reflects less variability in the regional response $\Theta_r$ (= Δ(R-F)_r/$\Delta T_{A,r}$) than in $\lambda_r$ (Fig. 5),

187  as well as less variability within the regions.  We therefore conclude that interannual variability

188  has less of an impact on Θ than λ.  We show additional evidence for the superior precision of Θ

189  in the Appendix.

190  As far as accuracy goes, we can compare Θ in the ensemble over the historical period to Θ in

191  response to much larger warming.  The ensemble median of Θ from the historic period (Fig. 2),

192  -1.04±0.01 W/m²/K (5-95% confidence interval), is close to the value obtained from an analysis

193  of the first 150 years of an abrupt $4xCO_2$ run of the same model, Θ = -1.03±0.04 W/m²/K, as

194  well as Θ calculated from all 2600 years of this run, Θ = -1.00±0.01 W/m²/K (values from the

195  $4xCO_2$ runs are all obtained using the Gregory method (Gregory et al., 2004) using annual

8

Deleted: ;

Deleted: t

Deleted: average

Deleted: average

Deleted: 4

201 average R and temperatures).  On the other hand, $\lambda$ changes substantially in the $4\times CO_2$ run as

202 the climate warms: $\lambda$ = -1.36$\pm$0.07 W/m$^2$/K when calculated from the first 150 years, but $\lambda$ =

203 -0.95$\pm$0.01 W/m$^2$/K from all 2600 years of that run.

204 We can verify this result in the CMIP5 abrupt $4\times CO_2$ ensemble.  It has been previously

205 demonstrated that plots of R-F vs. $T_S$ do not trace straight lines as the climate warms (Andrews

206 et al., 2015; Rugenstein et al., 2016; Rose and Rayborn, 2016; Armour, 2017), so $\lambda$ and ECS

207 calculated in a single model run may depend on the portion of the run selected.  In the CMIP5

208 abrupt $4\times CO_2$ ensemble, for example, average $\lambda$ calculated by regressing years 10-30 ($\lambda_{10\text{-}30}$) is

209 more negative than $\lambda$ calculated from years 30-150 ($\lambda_{30\text{-}150}$) by 0.49 W/m$^2$/K (Fig. 6).

210 Several explanations for this have been advanced, most prominently that $\lambda$ is function of the

211 pattern of surface warming (Senior and Mitchell, 2000; Armour et al., 2013; Andrews et al.,

212 2015; Gregory and Andrews, 2016; Zhou et al., 2016; Stevens et al., 2016).  Using $\Theta$ largely

213 eliminates this pattern effect: $\Theta_{10\text{-}30}$ and $\Theta_{30\text{-}150}$ have an average difference of 0.13 W/m$^2$/K for

214 the CMIP5 ensemble (Fig. 6).  Thus, we find additional evidence that $\Theta$ tends to be similar for

215 different amounts and patterns of warming.

216 The lack of curvature in the $\Theta$ calculations means there is curvature in the relation between $T_A$

217 and $T_S$ in the models.  Thus, the pattern effect's impact on ECS calculations shifts from $\lambda$ in the

218 traditional framework to the $\Delta T_S/\Delta T_A$ term in Eq. 4.  This also emphasizes the need to improve

219 our understanding of the factors that control $\Delta T_S/\Delta T_A$, as well as how future patterns of surface

220 warming will evolve.

221 There are several plausible reasons why $T_A$ may control R better than $T_S$.  It seems likely that

222 several of the feedbacks — e.g., lapse rate, water vapor, longwave cloud — should be more

223 strongly influenced by atmospheric temperatures than $T_S$.  More recently, it has been shown

224 that atmospheric temperatures also play a key role in regulating low clouds (Zhou et al., 2016,

225 2017), thereby influencing the shortwave cloud feedback.  This is also consistent with Ceppi et

226 al. (2017), who identified a dependence of ECS on atmospheric stability in models.  We have

227 not further investigated this — ultimately, our use of $T_A$ in Eq. 4 is based on observations

Deleted: .

Deleted: 50

Deleted: 6

231  (Murphy, 2010; Spencer and Braswell, 2010; Trenberth et al., 2015) that it correlates well with

232  R. Other metrics, such as global average atmospheric temperature work almost as well.

233  Clearly, further investigations on how to best describe the Earth's energy balance are

234  warranted.

235  Finally, one of our ultimate goals for this revised framework is to help produce better estimates

236  of ECS. We are working on a detailed analysis of ECS based on this framework and will publish

237  that in a follow-on paper, but we briefly show here how the advantages of the revised energy

238  balance framework may be leveraged to do this. Fig. 7a shows $\Theta$ calculated from control runs

239  of 25 CMIP5 models. To calculate $\Theta$ in the control runs, we break each control run into 17-year

240  segments and calculate monthly anomalies of $\Delta R$ and $\Delta T_A$ during each segment. Then, we

241  calculate $\Theta$ for each segment as the slope of the regression of $\Delta R$ vs. $\Delta T_A$ for that segment.

242  Thus, for each control run, we generate a large number of estimates of $\Theta$. The value in Fig. 7a is

243  the average of these individual values.

244  Fig. 7b shows the ECS of these models, calculated from the first 150 years of the abrupt 4x$CO_2$

245  runs using the Gregory method. If we assume that models with more accurate simulation of

246  short-term $\Theta$ produce more accurate estimates of ECS (Brown and Caldeira, 2017; Wu and

247  North, 2002), then we can use Figs. 7a and 7b to constrain ECS. We find that the 15 models

248  whose average short-term $\Theta$ falls within the uncertainty of $\Theta$ estimated from CERES

249  observations have ECS values ranging from 2.0-3.9 K, with an average of 2.9 K. This excludes

250  many of the highest ECS models, a result consistent with other analyses (Cox et al., 2018; Lewis

251  and Curry, 2015).

252  It would not have been possible to draw this conclusion with the conventional energy balance

253  framework. Fig. 7c shows the comparison between $\lambda$ from the control runs (calculated the

254  same way $\Theta$ was calculated) and CERES observations. Because of the much larger uncertainty

255  in the observational estimate of short-term $\lambda$, almost all models fall within the observational

256  range, thereby prohibiting any constraint on the ECS range.

**Deleted:** 6

**Deleted:** (Gregory et al., 2004)

**Deleted:** agrees with the

260  It may also be possible to use the relation between short-term and long-term $\Theta$ as an emergent

261  constraint to convert short-term observations to the long-term response.  There is some scatter

262  in the relation in the CMIP5 ensemble, however, so more analysis of how these relate is likely

263  required before ECS can be constrained in this way.

264  **Conclusions**

265  We have estimated ECS in each of a 100-member climate model ensemble using the same

266  energy-balance constraint used by many investigators to estimate ECS from 20[th]-century

267  historical observations.  We find that the method is imprecise — the estimates of ECS range

268  from 2.1 to 3.9 K (Fig. 2), with some ensemble members far from the model's true value of 2.9

269  K.  Given that we only have a single ensemble of reality, one should recognize that estimates of

270  ECS derived from the historical record may not be a good estimate of our climate system's true

271  value.

272  The source of the imprecision relates to the construction of the traditional energy balance

273  equation (Eq. 1).  In it, the response of TOA net flux (R-F) is parameterized in terms of global

274  average surface temperature ($T_S$).  Recent research has suggested that the response is not just

275  determined by the magnitude of $T_S$, but includes other factors, such as the pattern of $T_S$ (e.g.,

276  Armour et al., 2013; Andrews et al., 2015; Gregory and Andrews, 2016; Zhou et al., 2017) or the

277  lapse rate (e.g., Zhou et al., 2017; Ceppi and Gregory, 2017; Andrews and Webb, 2018).  As a

278  result, two ensemble members with the same $\Delta T_S$ can have different climate responses, $\Delta$(R-

279  F), leading to spread in the inferred $\lambda$.

280  The lack of a direct relationship between $T_S$ and radiation balance suggests that it may be

281  profitable to investigate alternative formulations. We test parameterizing the response in terms

282  of 500-hPa tropical temperature (Eq. 4) and find that it is superior in many ways.  Ultimately,

283  how investigators describe the energy balance of the planet will depend on the problem and

284  the available data.  The surface temperature is indeed special, so the traditional framework

285  may be preferred for some problems.  But investigators may find that the alternatives are

286  superior for certain problems, for instance constraining Earth's climate sensitivity.

**Deleted:** this suggests that some skepticism is appropriate when considering

289 _Appendix_
290 It has been previously noted in analyses of the historical record that $\lambda$ exhibits significant
291 interdecadal variability (Andrews et al., 2015; Gregory and Andrews, 2016; Zhou et al., 2016).
292 We can reproduce this in a 2000-year control run (a run with fixed pre-industrial boundary
293 conditions) of the MPI-ESM1.1 model. Fig. 8 shows $\lambda$ calculated in a sliding 17-year window

| Deleted: 6 |

294 and confirms significant temporal variability in $\lambda$. We can similarly calculate $\Theta$ and find that
295 temporal variability in $\Theta$ is substantially smaller (Fig. 8).

296 This result is reproduced in the CMIP5 control models. Fig. 9 plots the standard deviation of
297 each CMIP5 model's set of short-term $\lambda$ divided by the standard deviation of that model's set of
298 short-term $\Theta$ (as described previously, we calculate time series of short-term $\lambda$ and $\Theta$ values for
299 each model by regressing anomalies in a 17-year sliding window of the control runs). All of the

| Deleted: 6 |

300 models fall above 1, demonstrating that there is less variability in the $\Theta$ time series than in the
301 $\lambda$ time series in every climate model. This confirms that $\Theta$ is more robust with respect to
302 internal variability than $\lambda$. It also suggests that $\Theta$ estimated from the satellite data (Fig. 3)
303 should be considered a better estimate of the climate system's long-term value than $\lambda$
304 estimated from the same data set.

305

| Deleted: ▨▨▨▨▨▨▨▨▨▨▨▨▨▨▨▨▨Page Break▨▨▨▨▨▨▨▨▨▨▨▨▨▨▨▨ |

320 **References**

321 Aldrin, M., Holden, M., Guttorp, P., Skeie, R. B., Myhre, G., and Berntsen, T. K.: Bayesian
322     estimation of climate sensitivity based on a simple climate model fitted to observations of
323     hemispheric temperatures and global ocean heat content, Environmetrics, 23, 253-271,
324     10.1002/env.2140, 2012.
325 Andrews, T., Gregory, J. M., and Webb, M. J.: The dependence of radiative forcing and feedback
326     on evolving patterns of surface temperature change in climate models, J. Climate, 28,
327     1630-1648, 10.1175/JCLI-D-14-00545.1, 2015.
328 Andrews, T., and Webb, M. J.: The Dependence of Global Cloud and Lapse Rate Feedbacks on
329     the Spatial Structure of Tropical Pacific Warming, J. Climate, 31, 641-654, 10.1175/jcli-d-17-
330     0087.1, 2018.
331 Annan, J. D., and Hargreaves, J. C.: Using multiple observationally-based constraints to estimate
332     climate sensitivity, Geophys. Res. Lett., 33, 10.1029/2005gl025259, 2006.
333 Armour, K. C., Bitz, C. M., and Roe, G. H.: Time-varying climate sensitivity from regional
334     feedbacks, J. Climate, 26, 4518-4534, 10.1175/jcli-d-12-00544.1, 2013.
335 Armour, K. C.: Energy budget constraints on climate sensitivity in light of inconstant climate
336     feedbacks, Nature Clim. Change, 7, 331-335, 10.1038/nclimate3278, 2017.
337 Brown, P. T., and Caldeira, K.: Greater future global warming inferred from Earth's recent
338     energy budget, Nature, 552, 10.1038/nature24672, 2017.
339 Ceppi, P., and Gregory, J. M.: Relationship of tropospheric stability to climate sensitivity and
340     Earth's observed radiation budget, Proc. Natl. Acad. Sci., 10.1073/pnas.1714308114, 2017.
341 Collins, M., et al.: Long-term climate change: Projections, commitments and irreversibility, in:
342     Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the
343     Fifth Assessment Report of the Intergovernmental Panel on Climate Change, edited by:
344     Stocker, T. F., Qin, D., Plattner, G.-K., Tignor, M., Allen, S. K., Boschung, J., Nauels, A., Xia,
345     Y., Bex, V., and Midgley, P. M., Cambridge University Press, Cambridge, United Kingdom
346     and New York, NY, USA., 2013.
347 Cox, P. M., Huntingford, C., and Williamson, M. S.: Emergent constraint on equilibrium climate
348     sensitivity from global temperature variability, Nature, 553, 319-322,
349     10.1038/nature25450, 2018.
350 Dee, D. P., et al.: The ERA-Interim reanalysis: Configuration and performance of the data
351     assimilation system, Q. J. R. Meteor. Soc., 137, 553-597, 10.1002/qj.828, 2011.
352 Dessler, A. E., and Zelinka, M. D.: Climate feedbacks, in: Encyclopedia of Atmospheric Sciences,
353     edited by: North, G. R., Pyle, J., and Zhang, F., Elsevier, 18–25, 2014.
354 Forster, P. M.: Inference of climate sensitivity from analysis of Earth's energy budget, Annual
355     Review of Earth and Planetary Sciences, 44, 85-106, 10.1146/annurev-earth-060614-
356     105156, 2016.
357 Forster, P. M., et al.: Recommendations for diagnosing effective radiative forcing from climate
358     models for CMIP6, J. Geophys. Res., 121, 12460-12475, 10.1002/2016jd025320, 2016.
359 Giorgetta, M. A., et al.: Climate and carbon cycle changes from 1850 to 2100 in MPI-ESM
360     simulations for the Coupled Model Intercomparison Project phase 5, Journal of Advances in
361     Modeling Earth Systems, 5, 572-597, 10.1002/jame.20038, 2013.

362  Gregory, J. M., Stouffer, R. J., Raper, S. C. B., Stott, P. A., and Rayner, N. A.: An observationally
363      based estimate of the climate sensitivity, J. Climate, 15, 3117-3121, 10.1175/1520-
364      0442(2002)015<3117:aobeot>2.0.co;2, 2002.
365  Gregory, J. M., et al.: A new method for diagnosing radiative forcing and climate sensitivity,
366      Geophys. Res. Lett., 31, 10.1029/2003gl018747, 2004.
367  Gregory, J. M., and Andrews, T.: Variation in climate sensitivity and feedback parameters during
368      the historical period, Geophys. Res. Lett., 43, 3911-3920, 10.1002/2016GL068406, 2016.
369  Hansen, J., et al.: Efficacy of climate forcings, Journal of Geophysical Research: Atmospheres,
370      110, n/a-n/a, 10.1029/2005JD005776, 2005.
371  Hansen, J., Ruedy, R., Sato, M., and Lo, K.: Global surface temperature change, Rev. Geophys.,
372      48, 10.1029/2010rg000345, 2010.
373  Huber, M., Beyerle, U., and Knutti, R.: Estimating climate sensitivity and future temperature in
374      the presence of natural climate variability, Geophys. Res. Lett., 41, 2086-2092,
375      10.1002/2013GL058532, 2014.
376  Kummer, J. R., and Dessler, A. E.: The impact of forcing efficacy on the equilibrium climate
377      sensitivity, Geophys. Res. Lett., 41, 3565-3568, 10.1002/2014gl060046, 2014.
378  Lewis, N., and Curry, J. A.: The implications for climate sensitivity of AR5 forcing and heat
379      uptake estimates, Climate Dynamics, 45, 1009-1023, 10.1007/s00382-014-2342-y, 2015.
380  Loeb, N. G., et al.: Toward optimal closure of the Earth's top-of-atmosphere radiation budget, J.
381      Climate, 22, 748-766, 10.1175/2008jcli2637.1, 2009.
382  Mauritsen, T., et al.: Developments in the MPI-M Earth System Model version 1.2 (MPI-
383      ESM1.2), in preparation, 2018.
384  Murphy, D. M.: Constraining climate sensitivity with linear fits to outgoing radiation, Geophys.
385      Res. Lett., 37, 10.1029/2010GL042911, 2010.
386  Otto, A., et al.: Energy budget constraints on climate response, Nature Geoscience, 6, 415-416,
387      10.1038/ngeo1836, 2013.
388  Richardson, M., Cowtan, K., Hawkins, E., and Stolpe, M. B.: Reconciled climate response
389      estimates from climate models and the energy budget of Earth, Nature Clim. Change, 6,
390      931-935, 10.1038/nclimate3066, 2016.
391  Rose, B. E. J., and Rayborn, L.: The effects of ocean heat uptake on transient climate sensitivity,
392      Current Climate Change Reports, 2, 190-201, 10.1007/s40641-016-0048-4, 2016.
393  Rugenstein, M. A. A., Caldeira, K., and Knutti, R.: Dependence of global radiative feedbacks on
394      evolving patterns of surface heat fluxes, Geophys. Res. Lett., 43, 9877-9885,
395      10.1002/2016GL070907, 2016.
396  Santer, B. D., et al.: Statistical significance of trends and trend differences in layer-average
397      atmospheric temperature time series, J. Geophys. Res., 105, 7337-7356,
398      10.1029/1999jd901105, 2000.
399  Senior, C. A., and Mitchell, J. F. B.: The time-dependence of climate sensitivity, Geophys. Res.
400      Lett., 27, 2685-2688, 10.1029/2000GL011373, 2000.
401  Shindell, D. T.: Inhomogeneous forcing and transient climate sensitivity, 4, 274,
402      10.1038/nclimate2136, 2014.
403  Skeie, R. B., Berntsen, T., Aldrin, M., Holden, M., and Myhre, G.: A lower and more constrained
404      estimate of climate sensitivity using updated observations and detailed radiative forcing
405      time series, Earth System Dynamics, 5, 139-175, 10.5194/esd-5-139-2014, 2014.

406 Spencer, R. W., and Braswell, W. D.: On the diagnosis of radiative feedback in the presence of
407     unknown radiative forcing, J. Geophys. Res., 115, 10.1029/2009JD013371, 2010.
408 Stevens, B., Sherwood, S. C., Bony, S., and Webb, M. J.: Prospects for narrowing bounds on
409     Earth's equilibrium climate sensitivity, Earth's Future, 4, 512-522, 10.1002/2016EF000376,
410     2016.
411 Taylor, K. E., Stouffer, R. J., and Meehl, G. A.: An overview of CMIP5 and the experiment design,
412     Bull. Am. Met. Soc., 93, 485-498, 10.1175/bams-d-11-00094.1, 2012.
413 Trenberth, K. E., Zhang, Y., Fasullo, J. T., and Taguchi, S.: Climate variability and relationships
414     between top-of-atmosphere radiation and temperatures on Earth, J. Geophys. Res.,
415     10.1002/2014JD022887, 2015.
416 Wu, Q., and North, G. R.: Climate sensitivity and thermal inertia, Geophys. Res. Lett., 29,
417     10.1029/2002GL014864, 2002.
418 Zhou, C., Zelinka, M. D., and Klein, S. A.: Impact of decadal cloud variations on the Earth/'s
419     energy budget, Nature Geosci, 9, 871-874, 10.1038/ngeo2828, 2016.
420 Zhou, C., Zelinka, M. D., and Klein, S. A.: Analyzing the dependence of global cloud feedback on
421     the spatial pattern of sea surface temperature change with a Green's function approach,
422     Journal of Advances in Modeling Earth Systems, 9, 2174-2189, 10.1002/2017MS001096,
423     2017.
424
425

426

Fig. 1. Plot of annual and global average surface temperature from the 100 members of the
MPI-ESM1.1 ensemble (colored lines), along with the GISTEMP measurements (Hansen et
al., 2010) (white line).  Temperatures are referenced to the 1951-1980 average.

430



431

Figure 2. PDFs of (a) $\lambda$ (lighter) and $\Theta$ (darker) and (b) ECS derived from the members of the
MPI-ESM1.1 historical ensemble. The vertical lines are the 5[th], 50[th], and 95[th] percentile of each
distribution.

16

435



436 Figure 3. Scatter plot of detrended monthly anomalies of ΔR vs. (a) global average surface

437 temperature ΔT$_S$, (b) tropical average 500-hPa temperature ΔT$_A$. Observations cover the period

438 March 2000-July 2017 and anomalies are deviations from the mean annual cycle.  The dashed

439 lines are ordinary least-squares fits; the slope, 5-95% confidence interval, and correlation

440 coefficient are shown on each panel.  Confidence intervals account for autocorrelation of the

441 time series (Santer et al., 2000).

442
443

**Deleted:** Jan.

17

445

Fig. 4. Correlation coefficients between ΔR and temperature in CMIP5 control runs: black and red symbols represent the correlation with $\Delta T_S$ and $\Delta T_A$, respectively. The dot is the average of the correlation coefficients from the 17-year segments of the model run; the bars indicate the maximum and minimum values from the control run. The dashed lines are the corresponding correlation coefficients from the CERES regressions in Fig. 2.

451

459



460

Fig. 5. $\lambda_r$ and $\Theta_r$ calculated as regional average $\Delta$(R-F) divided by regional average temperature ($\Delta T_S$ for $\lambda$ and $\Delta T_A$ for $\Theta$). The regions are 90°S-19.4°S (SH), 19.4°S-19.4°N (EQ), and 19.4°N-90°N (NH). The values are calculated for each member of the 100-member ensemble; the solid symbols are the ensemble average while the bars show the 5-95% range.
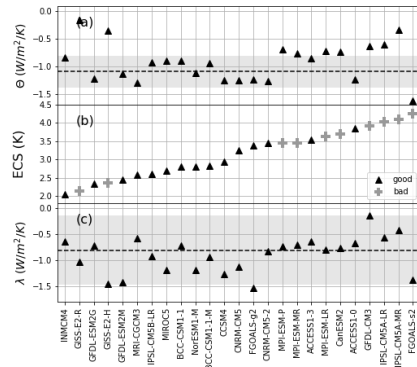


465

Fig. 6. Scatterplot of $\lambda_{10-30}$ vs. $\lambda_{30-150}$ (red circles) in CMIP5 abrupt4xCO$_2$ runs, as well as $\Theta_{10-30}$ vs. $\Theta_{30-150}$ (black triangles) in the same models. Each point represents one model. The dotted line is the 1:1 line. The subscripts (10-30, 30-150) indicate the years of the run from which the values are calculated.
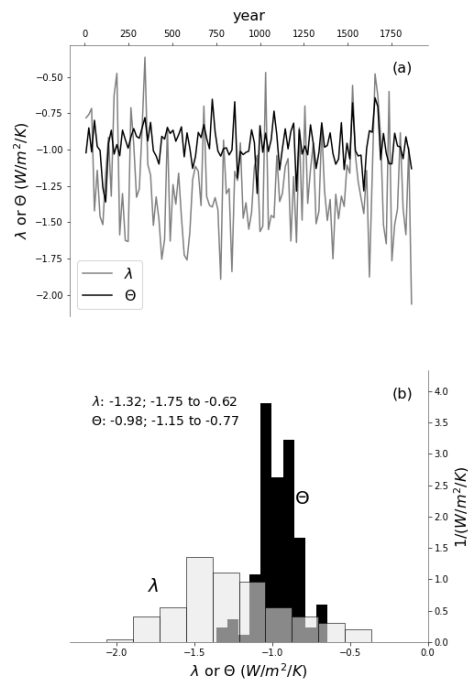
Deleted: gray

Deleted: crosses

472

Figure 7. (a) $\Theta$ from individual CMIP5 control runs. The dotted line is the estimate from

CERES observations; the gray region is the 5-95% confidence band. (b) ECS from each

CMIP5 model, estimated from the first 150 years of abrupt 4xCO$_2$ runs using the Gregory

method (Gregory et al., 2004). "Good" models are those whose $\Theta$ agrees with observations

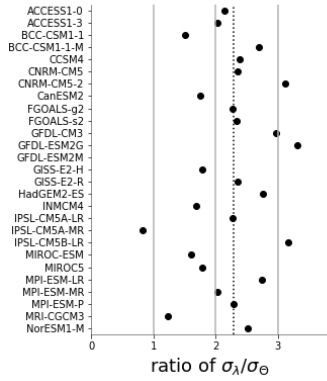in panel (a), "bad" models are those that do not. (c) Same as panel (a), but for $\lambda$.

478

Fig. 8. (a) Time series of λ (gray) and Θ (black) estimated in a 17-year sliding window of a 2000-year control run of the MPI-ESM1.1. (b) PDFs of the time series in panel a. Median and 5-95% confidence interval for each distribution is displayed on the plot.

485

486     Fig. 9. The standard deviation of the λ time series divided by the standard deviation of the

487     Θ time series.  Each time series is calculated from 17-year segments of CMIP5 control runs.

488     The dotted line is the ensemble average.

489