# The varying climate feedback parameter connected to the atmosphere-ocean coupling and ocean circulation

Diego Jiménez-de-la-Cuesta*

*Max-Planck-Institut für Meteorologie, Hamburg, Deutschland*

*Corresponding author*: Diego Jiménez-de-la-Cuesta, diego.jimenez@mpimet.mpg.de This work
has been submitted to the *Journal of Climate*. Copyright in this work may be transferred without
further notice.

## ABSTRACT

Models indicate a time-varying radiative response of the Earth system to $CO_2$ forcing (Andrews et al. 2012; Zhou et al. 2016). This variation implies a significant uncertainty in the estimates of climate sensitivity to increasing atmospheric $CO_2$ concentration (Hawkins and Sutton 2009; Grose et al. 2018). In energy-balance models, the temporal variation is represented as an additional feedback mechanism (Winton et al. 2010; Geoffroy et al. 2013a; Rohrschneider et al. 2019), which also depends on the ocean temperature change. Models and observations also indicate that a spatio-temporal pattern in surface warming controls this additional contribution to the radiative response by modulating the tropospheric instability (Ceppi and Gregory 2017; Zhou et al. 2016). Some authors focus on the atmospheric mechanisms that drive the feedback change (Stevens et al. 2016), reducing the role of the ocean's energy uptake variations. For the first time, I derive, using a linearized conceptual energy-balance model (Winton et al. 2010; Geoffroy et al. 2013a; Rohrschneider et al. 2019), an explicit mathematical expression of the radiative response and its temporal evolution. This expression connects the spatio-temporal warming pattern to a dynamical thermal capacity, stemming from changes in the ocean energy uptake. In comparison with more realistic energy-balance frameworks, and unlike the notion of additional feedback mechanisms, I show that an expanded effective thermal capacity better explains the variation of the radiative response, naturally connects with the spatio-temporal surface warming pattern and provides a non-circular framework to explain the variation of the climate feedback parameter.

26 *Significance statement.* Understanding the factors that change the Earth's radiative response to

27 forcing is central to reduce the uncertainty in the climate sensitivity estimates. The current

28 atmospheric-only view on the problem of the time-varying climate feedback parameter unneces-

29 sarily hides the ocean's role. This work shows a novel perspective for the problem, enabling the

30 development of a more general theory.

## 1. Introduction

31

32 Climate models show a wide range of temporal variation in their radiative response to $CO_2$ forcing

33 (Senior and Mitchell 2000; Andrews et al. 2012; Ceppi and Gregory 2017). This variation appears in

34 numerical experiments where the atmospheric $CO_2$ concentration is raised and maintained constant

35 afterwards. The rise in the atmospheric $CO_2$ concentration modifies the Earth's emissivity to long-

36 wave radiation, resulting in surface warming. Surface warming modifies the radiative flux at the

37 top of the atmosphere (TOA). The modified flux tends to cancel the energy imbalance introduced

38 by the radiative forcing. Surface warming also changes other variables, such as the atmospheric

39 temperature and humidity, that further modify the radiative flux. These changes are the feedback

40 mechanisms on surface warming. The net rate at which the globally-averaged surface warming

41 reduces the globally-averaged TOA imbalance is known as the climate feedback parameter.

42 If the feedback mechanisms did not change with time, the climate feedback parameter would be

43 constant, and a diagram of globally-averaged TOA imbalance change $N$ versus surface temperature

44 change $T_u$ ($NT-$diagram) would be linear. However, climate models present $NT-$diagrams with

45 different degrees of curvature, indicating a non-constant climate feedback parameter (presented in

46 Figure 1; Andrews et al. 2012; Ceppi and Gregory 2017). The degree of curvature is also modified

47 by forcing strength (Senior and Mitchell 2000; Meraner et al. 2013; Rohrschneider et al. 2019).

48 Temporal and state dependencies can be the root of climate feedback's variation. Observations

3

indicate that a spatio-temporal pattern of surface warming modifies cloud feedback in decadal timescales by altering atmospheric stability, leading to feedback changes that depend not only on surface warming (Zhou et al. 2016; Mauritsen 2016; Ceppi and Gregory 2017). The spatio-temporal warming pattern should depend also on the state of the system leading to contributions to the climate feedback's temporal and state dependencies.

In a globally-averaged energy-balance framework, $N$ should be equal to the forcing $F$ plus the radiative response of the system $R$, $N = F + R$. Following the classical picture of the linearized feedback mechanisms depending only on the surface warming (Gregory et al. 2002) $T_u$, we should have $R \sim \lambda T_u$, where $\lambda$ is a constant climate feedback parameter. Thus, if we consider a constant forcing $F$, the slope of the $NT-$diagram would be constant and equal to $\lambda$, in contradiction with observations and complex models as discussed above. Thus, either the non-linear component plays a more significant role, or the feedback mechanisms depend on more than the surface warming.

The problem cannot be solved by introducing more structure in the system. For instance, if we introduce two coupled layers (Winton et al. 2010; Geoffroy et al. 2013b; Rohrschneider et al. 2019). One layer represents the atmosphere, land and the mixed upper ocean layer: the upper layer. The deep layer corresponds to the deep ocean. Both layers have different thermal capacities. The thermal capacities provide two timescales. One can show that these timescales only delay the equilibrium but do not alter the climate feedback parameter $\lambda$. For the upper layer, the energy budget is $N_u = F + \lambda T_u - H$, where $H$ is the coupling between the upper and the deep layer: the deep-ocean energy uptake. The deep layer budget is, therefore, $N_d = H$. The sum of both budgets provides the planetary imbalance at the TOA $N = F + \lambda T_u$. Thus, $N$ has the same form as before.

Geoffroy et al. (2013a) introduced a perturbed deep-ocean energy uptake in the upper-layer budget, $H' := \varepsilon H$, where $\varepsilon$ is the efficacy parameter. The deep-layer budget is still equal to $H$, leading to a different energy imbalance at the TOA: $N = F + \lambda T_u + (H - H')$. Although the

4

73  $H - H'$ term appears to violate the energy conservation principle, it can be interpreted as additional

74  atmospheric feedbacks (e.g. Armour et al. 2013; Stevens et al. 2016) that depend on the deep

75  ocean's temperature. However, this explanation seems to be ad-hoc.

76  One of the weakness of the classical approach is the linear approximation. The climate feedback

77  parameter $\lambda$ is defined in terms of the derivative of $R$ evaluated at the initial state. In the same

78  manner, the deep-ocean energy uptake $H$ is defined in terms of a derivative evaluated at the initial

79  state. Therefore, the seemingly ad-hoc term introduced by Geoffroy et al. (2013a) has a more

80  concrete interpretation: the difference between the actual $H'$ and the reference $H$ results in a

81  surplus $(1 - \varepsilon)H$. This surplus comes from a dynamically-expanded capacity of the deep ocean

82  to uptake energy: an effective thermal capacity. This effective thermal capacity influences the

83  flux between the upper and deep layers, modifying the surface temperature from below. Once the

84  surface temperature is modified, the usual atmospheric feedbacks change.

85  Considering the analytical solutions of the modified two-layer model, I derive for the first time

86  an explicit mathematical expression for the slope of the $NT-$diagrams, including the explicit time

87  evolution of this slope. At its core, this expression has the ratio of change of the energy stored

88  by the upper and deep layers and supports the more concrete interpretation that I presented above.

89  The interpretation of these results is that the atmosphere-ocean coupling sets the spatio-temporal

90  warming pattern. Afterwards, the atmosphere adjusts, leading to the changes in the feedback

91  mechanisms. I first show the theory of this novel mathematical expression using the linearized

92  two-layer model. Afterwards, I take a step back away from the approximation to show why the

93  interpretation given above is more relevant.

5

## 2. Theory

The following equations define the modified linearized two-layer model (Geoffroy et al. 2013a)

$$
\begin{cases}
C_u\dot{T}_u = F + \lambda T_u - \varepsilon\gamma(T_u - T_d) \\
C_d\dot{T}_d = \qquad\qquad \gamma(T_u - T_d)
\end{cases}
$$

where the first equation corresponds to the upper-layer budget and the second equation to the deep layer. The parameters $\lambda$ and $\gamma$ are the climate feedback parameter and the rate of deep-ocean energy uptake in the neighbourhood of the initial state. $C_u$ and $C_d$ are respectively the thermal capacities (per unit area) of the upper and deep layers. $T_u$ and $T_d$ are temperature anomalies referred to the initial state and the dotted quantities are time total derivatives. The planetary imbalance is the sum of both equations, resulting in $N = F + \lambda T_u + (1 - \varepsilon)\gamma(T_u - T_d)$. Nonetheless, it is better to write these equations in the following fashion

$$
\begin{cases}
\dot{T}_u = F' + \lambda'T_u - \varepsilon\gamma'(T_u - T_d) \\
\dot{T}_d = \qquad\qquad \gamma'_d(T_u - T_d)
\end{cases}
\tag{1}
$$

where $F' := F/C_u$ with units of $\mathrm{K\,s^{-1}}$ and, $\lambda' := \lambda/C_u$, $\gamma' := \gamma/C_u$ and $\gamma'_d := \gamma/C_d$ with units of $\mathrm{s^{-1}}$. Equations (1) are a system of linear ordinary differential equations (Geoffroy et al. 2013a; Rohrschneider et al. 2019). Although the solutions are standard and widely discussed in other articles (e.g. Geoffroy et al. 2013a; Rohrschneider et al. 2019), here I will use the normal mode approach. In the following, I proceed by summarizing the relevant facts, leaving the full mathematical discussion to the appendix A of this article.

The homogeneous ($F' \equiv 0$) version of the system (1) has two distinct eigenvalues $\mu_\pm := (\hat{\lambda} \pm \kappa)/2$, where $\hat{\lambda} := \lambda' - \varepsilon\gamma' - \gamma'_d$ and $\kappa^2 := \hat{\lambda}^2 + 4\lambda'\gamma'_d$. These eigenvalues provide two distinct eigenvectors, forming a basis in which the full system (1) is uncoupled and, therefore, has a straight-forward solution. The eigensolutions $T_\pm$ are the solutions associated with each eigenvalue. Afterwards, one can return to the original representation, finding that $T_u$ and $T_d$ are linear combinations of

6

118  $T_\pm$. These linear combinations are the normal modes: the symmetric mode $T_s := T_+ + T_-$ and the

119  antisymmetric mode $T_a := T_+ - T_-$. The main result of this process is that $T_\mathrm{u} = T_s$ and $T_\mathrm{d} = \alpha T_s + \beta T_a$,

120  where $\alpha$ and $\beta$ are scalars that depend on the coefficients of the system (1). The normal-mode

121  representation makes explicit the coupling of the deep layer with the upper layer.

## 3. Results

123  *(i) The explicit slope of the NT−diagram*   From the solutions to system (1) written in terms of the

124  normal modes, one can obtain an expression for the slope of the $NT-$diagram, $\dot{N}/\dot{T}_\mathrm{u}$, of a system

125  under constant forcing. In the appendix, I derive the following closed expression for the slope

126  in terms of the derivatives of the normal modes and as a factor of the constant climate feedback

127  parameter $\lambda$

$$\frac{\dot{N}}{\dot{T}_\mathrm{u}} = \left\{ \frac{\varepsilon+1}{2\varepsilon} + \frac{\varepsilon-1}{2\varepsilon} \frac{C_\mathrm{u}\kappa}{|\lambda|} \left[ \left( \frac{\varepsilon}{C_\mathrm{u}} + \frac{1}{C_\mathrm{d}} \right) \frac{\gamma}{\kappa} - \frac{\dot{T}_a}{\dot{T}_s} \right] \right\} \lambda \qquad (2)$$

130  The main characteristic of equation (2) is the square-bracket term of its right-hand side. It contains

131  two parts. The first one sets a basic enhanced slope and contains the sum of the inverse of the

132  thermal capacities as if we had an electrical circuit with capacitors in series. The second part

133  provides the time evolution. It is a ratio of the changes in energy content. This ratio compares the

134  change in energy content of the deep layer with that of the upper layer. To confirm the importance

135  of the square-bracket term, one can take the limit as $\varepsilon \to 1$, where the pattern effect is cancelled in

136  equation (2)

$$\lim_{\varepsilon \to 1} \frac{\dot{N}}{\dot{T}_\mathrm{u}} = \lambda$$

139  The strong coupling between the upper and deep layers disappears. We end up with a constant

140  slope. However, if $\varepsilon \neq 1$, the climate feedback parameter varies with the ratio of the changes in

7

energy content from the deep to the upper layer around a basic value that depends on the thermal

capacities of the system, the square bracket term in equation (2).

*(ii) Explicit expression for the ratio term*    Using the full mathematical expressions for the solutions

$T_s$ and $T_a$ for constant forcing, I write the ratio term in equation (2) as

$$\frac{\dot{T}_a}{\dot{T}_s} = \tanh\left[\frac{\kappa}{2}(t - t_0) + \operatorname{arctanh}\left(\frac{\hat{\lambda} + 2\gamma'_d}{\kappa}\right)\right] \tag{3}$$

The ratio (3) grows in a sigmoidal fashion from $-1$ to $1$. This hyperbolic tangent has a scaling

factor $(\kappa/2)$ that sets the rate of change of the hyperbolic tangent between its extreme values. It also

has a shift (the arctanh term) that determines when the hyperbolic tangent crosses zero, governing

the contribution of the last term in equation (2). Both scaling and shift are in terms of $C_u$, $C_d$, $\varepsilon$, $\gamma$

and $\lambda$.

The interpretation of equation (3) is that, after the initial forcing, the deep ocean warms up slower

than the upper layer, steepening the slope of the $NT-$diagram. Once the ratio reaches the sign-

reversal point, the last term's contribution in equation (2) only flattens the slope of the $NT-$diagram.

The scaling factor and the shift of the ratio (3) set the timescale for the flattening. Equation (3)

expresses precisely the time evolution of the climate feedback parameter that others have only

approximated through numerical experiments with the modified two-layer model (Geoffroy et al.

2013a; Rohrschneider et al. 2019). Additionally, it establishes a third timescale in the Earth system,

related to the atmosphere-ocean coupling.

*(iii) Explicit expression of the climate feedback parameter*    Using the explicit expressions, the

equation (2) for the climate feedback parameter is

$$\frac{\dot{N}}{\dot{T}_u} = \frac{\varepsilon + 1}{2\varepsilon}\left(1 + \frac{\varepsilon - 1}{\varepsilon + 1}\frac{C_u\kappa}{|\lambda|}\left[\left(\frac{\varepsilon}{C_u} + \frac{1}{C_d}\right)\frac{\gamma}{\kappa} - \tanh\left(\frac{\kappa}{2}(t - t_0) + \operatorname{arctanh}\left(\frac{\hat{\lambda} + 2\gamma'_d}{\kappa}\right)\right)\right]\right)\lambda \tag{4}$$

8

164    The factor of $\lambda$ is composed of terms that are positive except for the ratio term coming from

165    equation (3). One can prove that at the start ($t = t_0$) the slope is

166
167
$$\frac{\dot{N}}{\dot{T}_\mathrm{u}}(t_0) = \left(1 + (\varepsilon - 1)\frac{\gamma}{|\lambda|}\right)\lambda$$

168    and from here up to the sign reversal of the ratio term, the slope flattens. The flattening is gentle

169    at first, but towards the sign reversal it accelerates.

170    At the time of sign reversal we have

171
172
$$\frac{\dot{N}}{\dot{T}_\mathrm{u}}(t_\mathrm{rev}) = \frac{\varepsilon + 1}{2\varepsilon}\left(1 + \frac{\varepsilon - 1}{\varepsilon + 1}\left(\frac{\varepsilon}{C_\mathrm{u}} + \frac{1}{C_\mathrm{d}}\right)\frac{C_\mathrm{u}\gamma}{|\lambda|}\right)\lambda$$

173    and from here and on, the ratio term becomes positive, leading to an even flatter slope. The flattening

174    decelerates and becomes gentle again. The asymptotic value of the slope of the $NT-$diagram is

175
176
$$\lim_{t\to\infty}\frac{\dot{N}}{\dot{T}_\mathrm{u}} = \frac{\varepsilon + 1}{2\varepsilon}\left(1 + \frac{\varepsilon - 1}{\varepsilon + 1}\frac{C_\mathrm{u}\kappa}{|\lambda|}\left[\left(\frac{\varepsilon}{C_\mathrm{u}} + \frac{1}{C_\mathrm{d}}\right)\frac{\gamma}{\kappa} - 1\right]\right)\lambda$$

177    *(iv) Numerical estimates of the atmosphere-ocean coupling*    By substituting in expression (4) the

178    parameter values found by Geoffroy et al. (2013a), I find the timescale for the sign reversal of the

179    $\dot{T}_a/\dot{T}_s$ ratio term. I use the multimodel mean values reported by Geoffroy et al. (2013a). For the

180    multimodel average values ($C_\mathrm{u} = 8.2\,\mathrm{W\,yr\,m^{-2}\,K^{-1}}$, $C_\mathrm{d} = 109\,\mathrm{W\,yr\,m^{-2}\,K^{-1}}$, $\gamma = 0.67\,\mathrm{W\,m^{-2}\,K^{-1}}$,

181    $\lambda = -1.18\,\mathrm{W\,m^{-2}\,K^{-1}}$ and $\varepsilon = 1.28$) the sign reversal of the ratio term takes place after 18.3 years.

182    This timescale lies between the fast (4.2 years) and slow (290 years) timescales established in terms

183    of the thermal capacities alone (Geoffroy et al. 2013a). This timescale is when we are at one half

184    of the change between the initial and final values of the slope.

185    I calculate the time for sign reversal using the rest of values in the tables of Geoffroy et al. (2013a)

186    and obtain that the multimodel average is 18.8 years. The minimum value is 8.8 years for GISS-

187    E2-R, whereas the maximum is 25.1 years for CNRM-CM5.1. If I compare with their estimates

188    of the fast and the slow timescales, even the extreme values fit well between both. Enlightening is

9

189 that the timescale of the sign reversal seems to fit with the de-facto 20-year standard to evaluate

190 the change in slope (e.g. Ceppi and Gregory 2017).

191 I also compare between the multimodel averages for all parameters and with the thermal ca-

192 pacities as calculated by Jiménez-de-la-Cuesta and Mauritsen (2019): $C_u = 7.2\,\mathrm{W\,yr\,m^{-2}\,K^{-1}}$,

193 $C_d = 367\,\mathrm{W\,yr\,m^{-2}\,K^{-1}}$. The calculated deep-layer thermal capacity is larger than the CMIP5

194 multi-model average, whereas the calculated value for the upper layer is smaller than the CMIP5

195 average. From these differences, we can note changes in the slope evolution (figure 2). Although

196 the difference in final slopes is small, the calculated thermal capacities strongly shift the sign-

197 reversal timescale: a deeper deep ocean lengthens the sign-reversal timescale, whereas a shallower

198 upper layer shortens it.

## 4. Analysis and Discussion

200 *(i) Consequences of the equation (4)*    We have two terms in the factor of equation (4): the identity

201 term and the $(\varepsilon - 1)$−term. The second term is only active if $\varepsilon \neq 1$, and has two contributions.

202 The first one is a constant contribution linked to the thermal capacities of the system. The second

203 contribution is time-varying and depends on the ratio $\dot{T}_a/\dot{T}_s$. This ratio measures the proportion of

204 energy that goes into the deep ocean compared to that stored in the upper layer. Together, these

205 terms provide a physical picture in which the slope's variation is determined by a basic thermal

206 capacity, which is expanded dynamically. The expansion stems from the changing energy fluxes

207 between the upper and deep layers that differ from the flux at the starting state represented by $\gamma$. This

208 flux difference changes the surface temperature and connects with the evolving spatio-temporal

209 warming pattern in a more concrete fashion, given that the evolution and spatial distribution of the

210 sea surface temperature corresponds to changes in the energy fluxes.

I showed above that the thermal capacities have a strong effect on the timescale at which the slope of the $NT-$diagram changes (figure 2). Thermal capacities in complex models depend strongly on the depth of the ocean mixed-layer and, therefore, on the atmosphere-ocean coupling, providing diverse behaviors (figure 3)

The above interpretation of the two-layer model in the context of the real Earth System is the following: The relative change in the energy fluxes due to the atmosphere-ocean coupling compels the atmospheric feedbacks to adjust. Thus, the magnitude of the changes in the atmospheric radiative response depend on the atmosphere-ocean coupling. One can argue that the non-local free-tropospheric warming from the deep-convective warm regions is enough to explain the change in feedbacks. However, for the pattern effect to act as proposed by Zhou et al. (2016); Mauritsen (2016); Ceppi and Gregory (2017), one needs the surface temperature variation in the subsidence areas. Here the atmosphere-ocean coupling and oceanic circulation enter to change the surface temperature, determining the spatio-temporal warming pattern. Thus, the atmosphere-ocean coupling can play a larger role than thought before (Kiehl 2007).

*(ii) State and forcing dependence* In this article, I ignored the dependence on the strength of forcing (Senior and Mitchell 2000; Meraner et al. 2013; Rohrschneider et al. 2019). However, such dependence should come from the reference values $\varepsilon$, $\lambda$ and $\gamma$ that exist under a particular forcing. Values of $\lambda$ and $\gamma$ are first-order derivatives in the neighbourhood of the starting states. Therefore, we need to explore the physics behind the $\varepsilon$ parameter to understand how it can change depending on forcing. The physics of $\varepsilon$ is the atmosphere-ocean coupling, including circulation. Therefore, we should explore how forcing impacts the atmosphere-ocean coupling, resulting in changes in the spatio-temporal warming pattern.

There are versions of linearized energy balance models in which a simple non-linear term is introduced (Rohrschneider et al. 2019). Although higher-order terms in the Taylor expansion of either the radiative response $R$ or the energy uptake $H$ can provide additional information, the dependencies arising from the atmosphere-ocean coupling, as shown in this article, are far more important in light of the results presented above.

*(iii) Non-linear planetary energy balance*   Above I presented evidence favouring the ocean's energy uptake central role in determining the spatio-temporal warming pattern and its effects on the atmospheric feedback mechanisms. I test this idea in a more general theoretical framework by writing the planetary energy budget in another widely-known representation

$$\frac{d}{dt}(CT_{\mathrm{u}}) = (1-\alpha)S + G - \epsilon\sigma(fT_{\mathrm{u}})^4 \tag{5}$$

where $S := S(t)$ in $\mathrm{W\,m^{-2}}$ is the incoming solar radiative flux at the TOA, $\alpha$ is the planetary albedo, $G := G(t)$ in $\mathrm{W\,m^{-2}}$ represents the remaining inputs (natural and anthropogenic), and the last term is the usual planetary long-wave emission, in $\mathrm{W\,m^{-2}}$, as a grey-body of emissivity $\epsilon$ and surface temperature $T_{\mathrm{u}}$ with $f$ the lapse-rate scaling factor for the emission temperature. At first inspection, we have the origin of the feedback mechanisms: the planetary albedo $\alpha$, the emissivity $\epsilon$ and the scaling factor $f$. On the one hand, we have the short-wave strand, the planetary albedo $\alpha := \alpha(T_{\mathrm{u}}, q_{cld,w}, \dots)$ that is a function of, e.g., the surface temperature (determining ice-sheet and sea-ice area) and the amount of liquid water in the atmosphere forming clouds. On the other hand, we have the long-wave thread, the emissivity and the lapse-rate scaling factor $\epsilon, f := \mathrm{f}(T_{\mathrm{u}}, q_v, q_{cld,w}, \dots)$, depending on, e.g. the surface temperature, and the amount of water vapor and cloud liquid water in the atmosphere.

The atmospheric feedback mechanisms cannot rely on any temperature we define inside the ocean. The ocean affects $\alpha$, $\epsilon$ and $f$ only through changing $T_{\mathrm{u}}$. In equation (5), we cannot see such

12

dependence. Therefore, here we would be tempted to artificially introduce it by saying that $\alpha$, $\epsilon$ and $f$ depend on another temperature in the ocean, as others have interpreted from the modified two-layer model. In this work, I have shown that there is another more natural place where the ocean enters into play: the energy imbalance at the TOA, $N$.

We usually picture $N$ as $N = (\mathrm{d}/\mathrm{d}t)(CT_u) = C\dot{T}_u$. However, we can have processes that modify the thermal capacity: For example, (a) if the mixed-layer depth varies, it modifies the amount of water in contact with the atmosphere; (b) deep-water upwelling not only reduces the surface temperature but increases the thermal capacity of the mixed layer; (c) a similar effect comes from water from ice sheet melting. Therefore, the sea surface temperature is controlled by the oceanic circulation and the ocean-atmosphere interaction, setting the warming pattern. Thus, a more proper definition for $N$ is $N = (\mathrm{d}/\mathrm{d}t)(CT_u) - \dot{C}T_u = C\dot{T}_u$. The term $\dot{C}T_u$ considers these processes that modify the sea surface temperature and the ocean energy fluxes from below.

If we rewrite equation (5) using the corrected $N$:

$$N = (1 - \alpha)S + G - \epsilon\sigma(fT_u)^4 - \dot{C}T_u \tag{6}$$

The last term of equation (6) is the representation of the effect of the spatio-temporal warming pattern. The factor $\dot{C}$ needs a new differential equation that describes the temporal variations of the energy uptake due to the ocean circulation and atmosphere-ocean interactions. If we want to use this non-linear framework, we would also need additional differential equations or constitutive relationships for $\alpha$, $\epsilon$ and $f$.

However, to look at the effect of this additional dynamical thermal capacity term on the atmospheric feedbacks, let us consider the total derivative of $N$. In the appendix B, I present the details

279 of this derivation. The expression for the total derivative of equation (6) divided by $\dot{T}_{\mathrm{u}}$ is

280
$$\frac{\dot{N}}{\dot{T}_{\mathrm{u}}} = \left[ (1 - \alpha) \frac{\dot{S}}{\dot{T}_{\mathrm{u}}} + \frac{\dot{G}}{\dot{T}_{\mathrm{u}}} \right]$$

281

282
$$- \left[ 1 + \frac{S}{4\epsilon\sigma f(fT_{\mathrm{u}})^3} \frac{\dot{\alpha}}{\dot{T}_{\mathrm{u}}} + \frac{T_{\mathrm{u}}}{4\epsilon} \frac{\dot{\epsilon}}{\dot{T}_{\mathrm{u}}} + \frac{T_{\mathrm{u}}}{f} \frac{\dot{f}}{\dot{T}_{\mathrm{u}}} + \frac{\dot{C}}{4\epsilon\sigma f(fT_{\mathrm{u}})^3} + \frac{T}{4\epsilon\sigma f(fT_{\mathrm{u}})^3} \frac{\ddot{C}}{\dot{T}_{\mathrm{u}}} \right] 4\epsilon\sigma f(fT_{\mathrm{u}})^3 \quad (7)$$

283 The first term is the contribution from the forcing variation, whereas the second term is the

284 contribution of the radiative response. This last term is the analogue for the expression of the

285 square bracket term in equation (2). This last term has a non-dimensional factor multiplying

286 the product of Planck feedback. In the non-dimensional factor we have a sum of terms. Each

287 term represents the contributions of feedbacks and other processes to the radiative response in

288 comparison to the Planck feedback. The first term in this factor is one, for the Planck feedback.

289 The second term is the contribution of the planetary-albedo feedback, including surface albedo as

290 well as short-wave cloud feedback. The third term is the contribution of the emissivity feedback,

291 including water-vapor feedback. The fourth term is the lapse-rate feedback. The last two terms

292 are the contribution of the atmosphere-ocean interaction and ocean circulation to the radiative

293 response.

294 Equation (7) is the non-linear analogue of equation (2) if we consider that the forcing is stationary.

295 In a linearization, the first four terms would provide a term analogous to $\lambda$. The last two terms,

296 considering that at least $\ddot{C} \neq 0$, would provide the analogue of the remaining terms corresponding

297 to a situation where $\varepsilon \neq 1$. If $\varepsilon = 1$, then $\dot{C} = 0$ and the linearization of equation (7), would

298 provide a constant feedback parameter if we accept that $\alpha$, $\epsilon$ and $f$ covary with $T_{\mathrm{u}}$ except for sign.

299 These facts provide more support for interpreting the modified two-layer model as introducing

300 a dynamical thermal capacity, turning the limelight to the role of the ocean circulation and the

301 atmosphere-ocean interactions in setting the sea surface temperature patterns.

14

## 5. Conclusions

I presented for the first time an explicit mathematical expression of the slope of the $NT-$diagrams. For that end, I used the linearized framework of the two-layer energy balance model. From the analysis, I uncover another timescale in the Earth System: the timescale of the changes of the climate feedback parameter. The timescale is related to the ratio of change in the energy content between the deep- and the upper-layer (atmosphere, land and mixed upper ocean). In CMIP5 models this time scale is around 18 years, providing theoretical support to the 20-year standard timescale used to study the change in the climate feedback parameter.

The mathematical expression and the analysis of the modified linearized two-layer model suggest that the spatio-temporal warming pattern is a product of the atmosphere-ocean coupling (the $H$ term in equations). Linearization uses the value of the coupling at the starting state. When introducing a modified coupling ($H'$) that represents departures from the starting state, the modified two-layer model can represent the variation of the slope of the $NT-$diagrams. The found mathematical expression suggest that the difference $H - H'$ acts as a dynamically-enlarged thermal capacity that depends on the ratio of change in the energy content between layers, changing the surface temperature. In the real Earth System, therefore, the atmosphere-ocean coupling and the circulation produce the evolving spatio-temporal warming pattern, to which the atmospheric feedbacks respond. I also shortly discussed this interpretation in terms of a non-linear framework, where feedback mechanisms not depending on the surface temperature are complicated to introduce. Instead of introducing exogenous dependencies in the feedback terms, the solution is that the thermal capacity of the system varies. The variation represents the atmosphere-ocean coupling and ocean circulation and relates to the spatio-temporal surface warming pattern. Its effect on the surface temperature then seamlessly translates to changes on the atmospheric feedback mechanisms.

15

*Data availability statement.* All the data used by this study is described fully in the main text. The numbers for reproducing the figure 2 and 3 are contained in the articles by Geoffroy et al. (2013a) and Jiménez-de-la-Cuesta and Mauritsen (2019). The theoretical considerations are fully described in the appendix to this article. Therefore, this article is fully reproducible and almost self contained.

## APPENDIX A

### Mathematical analysis of the modified two-layer model

In Classical Mechanics, a very coarse thinking would be reducing the field to the task of solving the equation $\dot{\mathbf{p}} = \mathbf{F}$ for any force term, either analytically or numerically. Going further leads to conservation principles and formulations of Classical Mechanics that provide more information without actually obtaining solutions, if that is possible at all. In this appendix, reduced to the scale of a simplified framework, I show that by delving deep into the mathematics of a system of linear ordinary differential equations, the structure of the solutions and its physical interpretation, one can obtain a new view on an old problem.

The appendix is written in an exhaustive way and I leave few things without development. The cases in which I do not show some algebraic step is because the necessary step has been already done or is very simple.

16

**Matrix form of the equations**

The equations of two-layer model Geoffroy et al. (2013a) are

$$N_u = C_u \dot{T}_u = F + \lambda T_u - \varepsilon\gamma(T_u - T_d)$$

$$N_d = C_d \dot{T}_d = \gamma(T_u - T_d) \tag{A1}$$

and the planetary imbalance is $N = N_u + N_d$. I present another form of the equations, where I divide by the thermal capacities.

$$\dot{T}_u = \frac{F}{C_u} + \frac{\lambda}{C_u}T_u - \varepsilon\frac{\gamma}{C_u}(T_u - T_d)$$

$$\dot{T}_d = \frac{\gamma}{C_d}(T_u - T_d)$$

If I define $F' := F/C_u$, $\lambda' := \lambda/C_u$, $\gamma' := \gamma/C_u$, $\gamma'_d := \gamma/C_d$, one can write the equations in a lean way

$$\dot{T}_u = F' + \lambda'T_u - \varepsilon\gamma'(T_u - T_d)$$

$$\dot{T}_d = \gamma'_d(T_u - T_d) \tag{A2}$$

I will put the system in matrix form. I define $\mathbf{T} := (T_u, T_d)$, $\mathbf{F}' := (F', 0)$ and

$$\mathbf{A} := \begin{pmatrix} \lambda' - \varepsilon\gamma' & \gamma'_d \\ \varepsilon\gamma' & -\gamma'_d \end{pmatrix} \tag{A3}$$

and the system can be written

$$\dot{\mathbf{T}} = \mathbf{F}' + \mathbf{TA} \tag{A4}$$

which is the representation of the system in the temperature basis.

**Eigenvalues and eigenvectors**

I want to analyse the normal modes of the system. For that end, I need the eigenvalues of the homogeneous system obtained as the solutions of the characteristic equation

$$(\lambda' - \varepsilon\gamma' - \mu)(-\gamma'_d - \mu) - \varepsilon\gamma'\gamma'_d = 0 \tag{A5}$$

17

369

$$-\lambda'\gamma'_d + \varepsilon\gamma'\gamma'_d + \mu\gamma'_d - \lambda'\mu + \varepsilon\gamma'\mu + \mu^2 - \varepsilon\gamma'\gamma'_d = 0$$

$$-\lambda'\gamma'_d + \mu\gamma'_d - \lambda'\mu + \varepsilon\gamma'\mu + \mu^2 = 0$$

$$-\lambda'\gamma'_d - (\lambda' - \varepsilon\gamma' - \gamma'_d)\mu + \mu^2 = 0$$

373

374 The solutions of equation (A5) are

$$\mu = \frac{(\lambda' - \varepsilon\gamma' - \gamma'_d) \pm \left[(\lambda' - \varepsilon\gamma' - \gamma'_d)^2 + 4\lambda'\gamma'_d\right]^{1/2}}{2} \tag{A6}$$

377 and, given that in the Earth $C_u < C_d$, one can prove that there are two real and different eigenval-

378 ues. One needs to check that the square root term is not complex or zero. This only happens if the

379 sum within the square root is negative or zero

$$(\lambda' - \varepsilon\gamma' - \gamma'_d)^2 + 4\lambda'\gamma'_d \le 0$$

$$(\lambda' - \varepsilon\gamma')^2 - 2(\lambda' - \varepsilon\gamma')\gamma'_d + \gamma'^2_d + 4\lambda'\gamma'_d \le 0$$

$$\lambda'^2 - 2\lambda'\varepsilon\gamma' + (\varepsilon\gamma')^2 - 2(\lambda' - \varepsilon\gamma')\gamma'_d + \gamma'^2_d + 4\lambda'\gamma'_d \le 0$$

$$\lambda'^2 - 2\lambda'\varepsilon\gamma' + (\varepsilon\gamma')^2 - 2\lambda'\gamma'_d + 2\varepsilon\gamma'\gamma'_d + \gamma'^2_d + 4\lambda'\gamma'_d \le 0$$

$$(\lambda'/\gamma'_d)^2 - 2(\lambda'/\gamma'_d)\varepsilon(\gamma'/\gamma'_d) + (\varepsilon(\gamma'/\gamma'_d))^2 + 2\varepsilon(\gamma'/\gamma'_d) + 1 + 2(\lambda'/\gamma'_d) \le 0$$

$$(\lambda'/\gamma'_d)^2 - 2(\lambda'/\gamma'_d)[\varepsilon(\gamma'/\gamma'_d) - 1] + (\varepsilon(\gamma'/\gamma'_d))^2 + 2\varepsilon(\gamma'/\gamma'_d) + 1 \le 0$$

$$(\lambda'/\gamma'_d)^2 - 2(\lambda'/\gamma'_d)[\varepsilon(\gamma'/\gamma'_d) - 1] + (\varepsilon(\gamma'/\gamma'_d) + 1)^2 \le 0$$

$$(\lambda'/\gamma'_d)^2 + (\varepsilon(C_d/C_u) + 1)^2 \le 2(\lambda'/\gamma'_d)[\varepsilon(C_d/C_u) - 1]$$

388

389  In the last inequality, the left-hand side is always positive. The right-hand side depends on the

390 sign of the factors. The middle factor is negative since $\lambda'$ is negative and $\gamma'_d$ is positive. The third

391 factor is positive provided that $\varepsilon > C_u/C_d$. Given that $\varepsilon \ge 1$ and $C_u < C_d$, then the third factor

392 is positive in our case. Then the right-hand side is negative. Thus, we obtained a contradiction

18

393 by supposing that the square root term was negative or zero. Therefore, the conclusion is that

394 the eigenvalues are two real and distinct numbers. Some CMIP5 models show $\varepsilon < 1$ according

395 to Geoffroy et al. (2013a). These also fit here. In the last condition of the above expression we

396 require that $\varepsilon(C_d/C_u) - 1 > 0$. If $\varepsilon \geq C_u/C_d$ this is fulfilled. $C_u/C_d$ is a small quantity and, in the

397 models that have a lesser than one $\varepsilon$, always the $\varepsilon$ is larger than this small quantity by an order of

398 magnitude. Thus, what I had said until now and will be said afterwards applies to all cases.

399 I call the solutions $\mu_+$ and $\mu_-$, depending on the sign of the square root term. Let us rewrite their

400 expression in more lean fashion. I define $\hat{\lambda} := \lambda' - \varepsilon\gamma' - \gamma'_d$ and we call $\kappa$ the square root term.

401 Then, I rewrite the solutions (A6) as

$$\mu_\pm = \frac{\hat{\lambda} \pm \kappa}{2} \tag{A7}$$

402
403

404 Now that I know the eigenvalues, one should get the eigenvectors of the system and solve it

405 easily. The eigenvectors are the generators of the kernel of the operators $\mathbf{A} - \mu_\pm \mathrm{id}$. Let us write

406 the diagonal of the matrix $\mathbf{A}$ with the definition of $\hat{\lambda}$

$$\mathbf{A} = \begin{pmatrix} \hat{\lambda} + \gamma'_d & \gamma'_d \\ \varepsilon\gamma' & \hat{\lambda} - (\lambda' - \varepsilon\gamma') \end{pmatrix}$$

407

408

409 and then the matrices for each eigenvalue have the form

$$\mathbf{A} - \mu_\pm \mathrm{id} = \begin{pmatrix} \hat{\lambda} + \gamma'_d - \mu_\pm & \gamma'_d \\ \varepsilon\gamma' & \hat{\lambda} - (\lambda' - \varepsilon\gamma') - \mu_\pm \end{pmatrix}$$

410

$$= \begin{pmatrix} \mu_\mp + \gamma'_d & \gamma'_d \\ \varepsilon\gamma' & \mu_\mp - (\lambda' - \varepsilon\gamma') \end{pmatrix}$$

411

412

413 Since eigenvalues are real and distinct, there should be two linearly-independent eigenvectors,

414 one for each eigenvalue. These vectors should fulfill that $\mathbf{e}_\pm(\mathbf{A} - \mu_\pm \mathrm{id}) = 0$. Solving that linear

19

system, I find the eigenvectors in temperature representation

$$\mathbf{e}_\pm = \mathbf{e}_u - \frac{\mu_\mp + \gamma'_d}{\varepsilon\gamma'}\mathbf{e}_d \tag{A8}$$

The procedure to get the result is to solve the system of homogeneous linear equations $\mathbf{e}_\pm(\mathbf{A} - \mu_\pm\,\mathrm{id}) = 0$

$$\begin{cases} (\mu_\mp + \gamma'_d)e_{\pm,u} & +\varepsilon\gamma' e_{\pm,d} = 0 \\[2mm] \gamma'_d e_{\pm,u} + [\mu_\mp - (\lambda' - \varepsilon\gamma')]e_{\pm,d} = 0 \end{cases}$$

I solve the first equation for the component $e_{\pm,d}$, and substitute this result on the second equation

$$e_{\pm,d} = -\frac{\mu_\mp + \gamma'_d}{\varepsilon\gamma'}e_{\pm,u} \longrightarrow$$

$$\left(\gamma'_d - \frac{[\mu_\mp - (\lambda' - \varepsilon\gamma')](\mu_\mp + \gamma'_d)}{\varepsilon\gamma'}\right)e_{\pm,u} = 0$$

$$\frac{\varepsilon\gamma'\gamma'_d - [\mu_\mp - (\lambda' - \varepsilon\gamma')](\mu_\mp + \gamma'_d)}{\varepsilon\gamma'}e_{\pm,u} = 0, (\varepsilon,\gamma' \neq 0)\ \therefore$$

$$\left\{\varepsilon\gamma'\gamma'_d - [\mu_\mp - (\lambda' - \varepsilon\gamma')](\mu_\mp + \gamma'_d)\right\}e_{\pm,u} = 0$$

$$\left\{\varepsilon\gamma'\gamma'_d + [(\lambda' - \varepsilon\gamma') - \mu_\mp](\gamma'_d + \mu_\mp)\right\}e_{\pm,u} = 0$$

$$-\left\{-\varepsilon\gamma'\gamma'_d + [(\lambda' - \varepsilon\gamma') - \mu_\mp](-\gamma'_d - \mu_\mp)\right\}e_{\pm,u} = 0$$

and in the last expression we have two options: either $e_{\pm,u}$ is zero or the term within curly braces is zero. However, the expression in curly braces is the characteristic equation (A5) and then always vanishes identically. This means that $e_{\pm,u} = \alpha \in \mathbb{R}$ can be chosen arbitrarily. I plug in this result in the expression for $e_{\pm,d}$ and get that

$$e_{\pm,u} = \alpha$$

$$e_{\pm,d} = -\frac{\mu_\mp + \gamma'_d}{\varepsilon\gamma'}\alpha$$

20

439    or as a vector in the temperature basis

440
$$\mathbf{e}_\pm = e_{\pm,u}\mathbf{e}_u + e_{\pm,d}\mathbf{e}_d$$

441
$$\mathbf{e}_\pm = \alpha\mathbf{e}_u - \frac{\mu_\mp + \gamma'_d}{\varepsilon\gamma'}\alpha\mathbf{e}_d$$
442

443    and since $\alpha$ is arbitrary this means we are in front of a subspace of vectors. I choose a basis by

444    selecting $\alpha = 1$.

445
$$\mathbf{e}_\pm = \mathbf{e}_u - \frac{\mu_\mp + \gamma'_d}{\varepsilon\gamma'}\mathbf{e}_d$$
446

447    which is the same as the equation (A8).

448    Now, I can derive the expressions of the temperature basis vectors in terms of the two eigenvectors.

449    If one solves for $e_u$ in equation (A8)

450
$$\mathbf{e}_\pm + \frac{\mu_\mp + \gamma'_d}{\varepsilon\gamma'}\mathbf{e}_d = \mathbf{e}_u$$
451

452    but we have here two expressions in a condensed way. Therefore,

453
$$\mathbf{e}_- + \frac{\mu_+ + \gamma'_d}{\varepsilon\gamma'}\mathbf{e}_d = \mathbf{e}_+ + \frac{\mu_- + \gamma'_d}{\varepsilon\gamma'}\mathbf{e}_d$$

454
$$\left(\frac{\mu_+ + \gamma'_d}{\varepsilon\gamma'} - \frac{\mu_- + \gamma'_d}{\varepsilon\gamma'}\right)\mathbf{e}_d = \mathbf{e}_+ - \mathbf{e}_-$$

455
$$\frac{(\mu_+ + \gamma'_d) - (\mu_- + \gamma'_d)}{\varepsilon\gamma'}\mathbf{e}_d = \mathbf{e}_+ - \mathbf{e}_-$$

456
$$\frac{\mu_+ - \mu_-}{\varepsilon\gamma'}\mathbf{e}_d = \mathbf{e}_+ - \mathbf{e}_-$$

457
$$\mathbf{e}_d = \frac{\varepsilon\gamma'}{\mu_+ - \mu_-}(\mathbf{e}_+ - \mathbf{e}_-)$$
458

459    Thus, I have expressed $\mathbf{e}_d$ in terms of the eigenvectors.

21

460    Now, I substitute the last result on one of the expressions for $\mathbf{e}_u$.

461
$$\mathbf{e}_+ + \frac{\mu_- + \gamma'_d}{\varepsilon\gamma'}\mathbf{e}_d = \mathbf{e}_u$$

462
$$\mathbf{e}_+ + \frac{\mu_- + \gamma'_d}{\varepsilon\gamma'}\frac{\varepsilon\gamma'}{\mu_+ - \mu_-}(\mathbf{e}_+ - \mathbf{e}_-) = \mathbf{e}_u$$

463
$$\mathbf{e}_+ + \frac{\mu_- + \gamma'_d}{\mu_+ - \mu_-}(\mathbf{e}_+ - \mathbf{e}_-) = \mathbf{e}_u$$

464
$$\left(1 + \frac{\mu_- + \gamma'_d}{\mu_+ - \mu_-}\right)\mathbf{e}_+ - \frac{\mu_- + \gamma'_d}{\mu_+ - \mu_-}\mathbf{e}_- = \mathbf{e}_u$$

465
$$\frac{\mu_+ - \mu_- + \mu_- + \gamma'_d}{\mu_+ - \mu_-}\mathbf{e}_+ - \frac{\mu_- + \gamma'_d}{\mu_+ - \mu_-}\mathbf{e}_- = \mathbf{e}_u$$

466
467
$$\frac{\mu_+ + \gamma'_d}{\mu_+ - \mu_-}\mathbf{e}_+ - \frac{\mu_- + \gamma'_d}{\mu_+ - \mu_-}\mathbf{e}_- = \mathbf{e}_u$$

468    and the temperature basis vectors in the eigenvector representation are

469
$$\mathbf{e}_u = \frac{\mu_+ + \gamma'_d}{\mu_+ - \mu_-}\mathbf{e}_+ - \frac{\mu_- + \gamma'_d}{\mu_+ - \mu_-}\mathbf{e}_-$$

(A9)

470
$$\mathbf{e}_d = \frac{\varepsilon\gamma'}{\mu_+ - \mu_-}(\mathbf{e}_+ - \mathbf{e}_-)$$

## 471    Matrix in the eigenvector representation. Solutions

472    With these results, I can write the matrix **A** (A3) in the eigenvector basis and it should be the

473    following diagonal matrix

474
$$\mathbf{B} = \begin{pmatrix} \mu_+ & 0 \\ 0 & \mu_- \end{pmatrix}$$

(A10)

475

476    I show how to get to this result. Let subscripts represent rows and superscripts represent columns.

477    I define that latin indices $(i, j, k, \dots)$ have the possible values u,d; and greek indices $(\alpha, \beta, \zeta \dots)$

478    have possible values $+, -$. Also, repeated indices in expressions mean summation over the set of

479    possible values. With these considerations, equation (A9) is

480
481
$$\mathbf{e}_i = \Lambda_i^\alpha \mathbf{e}_\alpha$$

22

where the rows of matrix $\Lambda$ contain the coordinates of each of the vectors of the temperature basis in the eigenvector representation. Analogously, equation (A8) is

$$\mathbf{e}_\alpha = \Theta_\alpha^i \mathbf{e}_i$$

where matrix $\Theta$ has in its rows the coordinates the eigenvector basis in the temperature representation. This means that

$$\mathbf{e}_\alpha = \Theta_\alpha^i \mathbf{e}_i = \Theta_\alpha^i \Lambda_i^\beta \mathbf{e}_\beta$$

which is only possible if the matrices $\Lambda$ and $\Theta$ are inverse of each other

$$\mathbf{e}_\alpha = \delta_\alpha^\beta \mathbf{e}_\beta = \mathbf{e}_\alpha$$

Thus, we write $\Theta = \Lambda^{-1}$.

Now, matrix **A** is the temperature representation of a linear operator $f$. If $\mathbf{v} = v^j \mathbf{e}_j$ is a vector in the temperature representation, then the action of the linear operator $f$ should be $f(\mathbf{v}) = f(v^j \mathbf{e}_j) = v^j f(\mathbf{e}_j)$. Then the action of $f$ on a vector expressed in a given basis only depends on the action of the operator on the basis: $f(\mathbf{v}) = f(v^j \mathbf{e}_j) = v^j f(\mathbf{e}_j) = v^j \mathbf{A}_j^k \mathbf{e}_k$. Thus, the matrix **A** has in its rows the coordinates in the temperature representation of the action of $f$ over each basis vector. Once one understands what is happening under the hood, what we want is the matrix **B**, which is the representation of $f$ in the eigenvector basis. Therefore, I begin with the basic relationship in the temperature representation and introduce the change of representation using the alternative

23

502 representation of equations (A8) and (A9)

$$f(\mathbf{e}_i) = \mathbf{A}_i^j \Lambda_j^\zeta \mathbf{e}_\zeta$$

503

$$f(\Lambda_i^\alpha \mathbf{e}_\alpha) = \mathbf{A}_i^j \Lambda_j^\zeta \mathbf{e}_\zeta$$

504

$$\Lambda_i^\alpha f(\mathbf{e}_\alpha) = \mathbf{A}_i^j \Lambda_j^\zeta \mathbf{e}_\zeta$$

505

$$(\Lambda^{-1})_\beta^i \Lambda_i^\alpha f(\mathbf{e}_\alpha) = (\Lambda^{-1})_\beta^i \mathbf{A}_i^j \Lambda_j^\zeta \mathbf{e}_\zeta$$

506

$$f(\mathbf{e}_\beta) = (\Lambda^{-1})_\beta^i \mathbf{A}_i^j \Lambda_j^\zeta \mathbf{e}_\zeta, \, f(\mathbf{e}_\beta) := \mathbf{B}_\beta^\zeta \mathbf{e}_\zeta$$

507

$$\mathbf{B}_\beta^\zeta = (\Lambda^{-1})_\beta^i \mathbf{A}_i^j \Lambda_j^\zeta$$

508
509

510 or in matrix notation $\mathbf{B} = \Lambda^{-1} \mathbf{A} \Lambda$. Then, I multiply the matrices

511
$$\Lambda^{-1} = \begin{pmatrix} 1 & -\frac{\mu_- + \gamma'_d}{\varepsilon\gamma'} \\ 1 & -\frac{\mu_+ + \gamma'_d}{\varepsilon\gamma'} \end{pmatrix}, \mathbf{A} = \begin{pmatrix} \hat{\lambda} + \gamma'_d & \gamma'_d \\ \varepsilon\gamma' & -\gamma'_d \end{pmatrix}, \Lambda = \begin{pmatrix} \frac{\mu_+ + \gamma'_d}{\mu_+ - \mu_-} & -\frac{\mu_- + \gamma'_d}{\mu_+ - \mu_-} \\ \frac{\varepsilon\gamma'}{\mu_+ - \mu_-} & -\frac{\varepsilon\gamma'}{\mu_+ - \mu_-} \end{pmatrix}$$

512

513 First, note that $\mu_+ - \mu_- = \kappa$. One also looks at the following quantities that will help in the

514 process: $\mu_+ + \mu_- = \hat{\lambda}$ and $\mu_+ \mu_- = \frac{1}{4}(\hat{\lambda}^2 - \kappa^2) = \frac{1}{4}(\hat{\lambda}^2 - \hat{\lambda}^2 - 4\lambda'\gamma'_d) = -\lambda'\gamma'_d$. I proceed with the first

515 product, $\Lambda^{-1}\mathbf{A}$.

516
$$\Lambda^{-1}\mathbf{A} = \begin{pmatrix} 1 & -\frac{\mu_- + \gamma'_d}{\varepsilon\gamma'} \\ 1 & -\frac{\mu_+ + \gamma'_d}{\varepsilon\gamma'} \end{pmatrix} \begin{pmatrix} \hat{\lambda} + \gamma'_d & \gamma'_d \\ \varepsilon\gamma' & -\gamma'_d \end{pmatrix}$$

517
$$= \begin{pmatrix} \hat{\lambda} + \gamma'_d - \mu_- - \gamma'_d & \left(1 + \frac{\mu_- + \gamma'_d}{\varepsilon\gamma'}\right)\gamma'_d \\ \hat{\lambda} + \gamma'_d - \mu_+ - \gamma'_d & \left(1 + \frac{\mu_+ + \gamma'_d}{\varepsilon\gamma'}\right)\gamma'_d \end{pmatrix}$$

518
$$= \begin{pmatrix} \hat{\lambda} - \mu_- & \frac{\varepsilon\gamma' + \mu_- + \gamma'_d}{\varepsilon\gamma'}\gamma'_d \\ \hat{\lambda} - \mu_+ & \frac{\varepsilon\gamma' + \mu_+ + \gamma'_d}{\varepsilon\gamma'}\gamma'_d \end{pmatrix}$$

519
$$= \begin{pmatrix} \mu_+ & \frac{\varepsilon\gamma' + \mu_- + \gamma'_d}{\varepsilon\gamma'}\gamma'_d \\ \mu_- & \frac{\varepsilon\gamma' + \mu_+ + \gamma'_d}{\varepsilon\gamma'}\gamma'_d \end{pmatrix}$$

520

24

521 and multiply the result by $\Lambda$

522
$$\Lambda^{-1}\mathbf{A}\Lambda = \begin{pmatrix} \mu_+ & \frac{\varepsilon\gamma'+\mu_-+\gamma'_d}{\varepsilon\gamma'}\gamma'_d \\ \mu_- & \frac{\varepsilon\gamma'+\mu_++\gamma'_d}{\varepsilon\gamma'}\gamma'_d \end{pmatrix}\begin{pmatrix} \frac{\mu_++\gamma'_d}{\mu_+-\mu_-} & -\frac{\mu_-+\gamma'_d}{\mu_+-\mu_-} \\ \frac{\varepsilon\gamma'}{\mu_+-\mu_-} & -\frac{\varepsilon\gamma'}{\mu_+-\mu_-} \end{pmatrix}$$

523
$$= \frac{1}{\kappa}\begin{pmatrix} \mu_+^2 + \mu_+\gamma'_d + \varepsilon\gamma'\gamma'_d + \mu_-\gamma'_d + \gamma'^2_d & -\mu_+\mu_- - \mu_+\gamma'_d - \varepsilon\gamma'\gamma'_d - \mu_-\gamma'_d - \gamma'^2_d \\ \mu_-\mu_+ + \mu_-\gamma'_d + \varepsilon\gamma'\gamma'_d + \mu_+\gamma'_d + \gamma'^2_d & -\mu_-^2 - \mu_-\gamma'_d - \varepsilon\gamma'\gamma'_d - \mu_+\gamma'_d - \gamma'^2_d \end{pmatrix}$$

524
$$= \frac{1}{\kappa}\begin{pmatrix} \mu_+^2 + (\hat{\lambda} + \varepsilon\gamma' + \gamma'_d)\gamma'_d & -\mu_+\mu_- - (\hat{\lambda} + \varepsilon\gamma' + \gamma'_d)\gamma'_d \\ \mu_-\mu_+ + (\hat{\lambda} + \varepsilon\gamma' + \gamma'_d)\gamma'_d & -\mu_-^2 - (\hat{\lambda} + \varepsilon\gamma' + \gamma'_d)\gamma'_d \end{pmatrix}$$

525
$$= \frac{1}{\kappa}\begin{pmatrix} \mu_+^2 - \mu_+\mu_- & \lambda'\gamma'_d - \lambda'\gamma'_d \\ -\lambda'\gamma'_d + \lambda'\gamma'_d & -\mu_-^2 + \mu_+\mu_- \end{pmatrix} = \frac{1}{\kappa}\begin{pmatrix} \mu_+\kappa & 0 \\ 0 & \mu_-\kappa \end{pmatrix} = \begin{pmatrix} \mu_+ & 0 \\ 0 & \mu_- \end{pmatrix}$$

526

527 the last line is the result that we wanted to check.

528 In the eigenvector representation the system (A4) has the following form

529
$$\dot{\mathbf{T}} = \mathbf{F}' + \mathbf{TB} \tag{A11}$$
530

531 and, therefore, is decoupled. Therefore, I can solve each equation separately. I only need to

532 transform the forcing vector to the eigenvector representation.

533 The equations are

534
$$\dot{T}_\pm = F'_\pm + \mu_\pm T_\pm$$
535

536 and the solutions of a generic initial value problem are

537
$$T_\pm = \left(T_{\pm,0} + \int_{t_0}^{t} F'_\pm e^{-\mu_\pm(\tau-t_0)}\mathrm{d}\tau\right)e^{\mu_\pm(t-t_0)} \tag{A12}$$
538

539 where the initial values in the eigenvector representation in terms of the initial values in the

540 temperature representation are

541
$$T_{\pm,0} = \pm\frac{1}{\mu_+ - \mu_-}[(\mu_\pm + \gamma'_d)T_{u,0} + \varepsilon\gamma' T_{d,0}]$$
542

25

the forcing components are

$$F'_\pm = \pm \frac{\mu_\pm + \gamma'_d}{\mu_+ - \mu_-} F'$$

and the solutions in the temperature representation are

$$T_u = T_+ + T_-$$

$$T_d = -\frac{\mu_- + \gamma'_d}{\varepsilon\gamma'} T_+ - \frac{\mu_+ + \gamma'_d}{\varepsilon\gamma'} T_-$$

If I further expand the $T_d$ solution, the form of the solutions is more elegant

$$T_u = T_+ + T_-$$

$$T_d = -\frac{\hat{\lambda} + 2\gamma'_d}{2\varepsilon\gamma'}(T_+ + T_-) + \frac{\kappa}{2\varepsilon\gamma'}(T_+ - T_-) \tag{A13}$$

since it shows that the solutions in the temperature space are in a sort of symmetric and antisymmetric combinations of the solutions in the eigenvector representation. These are the normal modes. One thing to note is that the upper temperature is the symmetric mode and the deep temperature is a mixture of symmetric and antisymmetric modes.

I show how I got the solutions (A13). Just expand the $T_d$ equation.

$$
\begin{aligned}
T_d &= -\frac{\mu_- + \gamma'_d}{\varepsilon\gamma'} T_+ - \frac{\mu_+ + \gamma'_d}{\varepsilon\gamma'} T_- \\
&= -\frac{1}{\varepsilon\gamma'}\left[\left(\frac{\hat{\lambda} - \kappa}{2} + \gamma'_d\right)T_+ + \left(\frac{\hat{\lambda} + \kappa}{2} + \gamma'_d\right)T_-\right] \\
&= -\frac{1}{\varepsilon\gamma'}\left[\left(\frac{\hat{\lambda} + 2\gamma'_d}{2} - \frac{\kappa}{2}\right)T_+ + \left(\frac{\hat{\lambda} + 2\gamma'_d}{2} + \frac{\kappa}{2}\right)T_-\right] \\
&= -\frac{1}{2\varepsilon\gamma'}\left[(\hat{\lambda} + 2\gamma'_d)(T_+ + T_-) - \kappa(T_+ - T_-)\right]
\end{aligned}
$$

From now on, I write $T_s := T_+ + T_-$ and $T_a := T_+ - T_-$.

26

**Planetary imbalance**

Now, I will find an expression for the planetary imbalance in terms of the equations (A13). The mathematical expression that I should expand is $N = N_{\mathrm{u}} + N_{\mathrm{d}} = C_{\mathrm{u}}\dot{T}_{\mathrm{u}} + C_{\mathrm{d}}\dot{T}_{\mathrm{d}}$

$$C_{\mathrm{u}}\dot{T}_{\mathrm{u}} = C_{\mathrm{u}}\dot{T}_s$$

$$C_{\mathrm{d}}\dot{T}_{\mathrm{d}} = -C_{\mathrm{d}}\frac{\hat{\lambda}+2\gamma'_d}{2\varepsilon\gamma'}\dot{T}_s + C_{\mathrm{d}}\frac{\kappa}{2\varepsilon\gamma'}\dot{T}_a \therefore$$

$$N = C_{\mathrm{u}}\dot{T}_s - C_{\mathrm{d}}\frac{\hat{\lambda}+2\gamma'_d}{2\varepsilon\gamma'}\dot{T}_s + C_{\mathrm{d}}\frac{\kappa}{2\varepsilon\gamma'}\dot{T}_a$$

$$= \left(C_{\mathrm{u}} - C_{\mathrm{d}}\frac{\hat{\lambda}+2\gamma'_d}{2\varepsilon\gamma'}\right)\dot{T}_s + C_{\mathrm{d}}\frac{\kappa}{2\varepsilon\gamma'}\dot{T}_a$$

$$= \mathrm{C}_s\dot{T}_s + \mathrm{C}_a\dot{T}_a$$

Now, $\dot{T}_{\pm} = F'_{\pm} + \mu_{\pm}T_{\pm}$, then

$$\dot{T}_s = \mu_+T_+ + \mu_-T_- + (F'_+ + F'_-) = \mu_+T_+ + (\mu_+ - \kappa)T_- + (F'_+ + F'_-)$$

$$= \mu_+T_s - \kappa T_- + (F'_+ + F'_-) = \mu_+T_s - \frac{\kappa}{2}(T_s - T_a) + (F'_+ + F'_-)$$

$$= \frac{\hat{\lambda}}{2}T_s + \frac{\kappa}{2}T_a + (F'_+ + F'_-) = \frac{\hat{\lambda}}{2}T_s + \frac{\kappa}{2}T_a + F'$$

$$\dot{T}_a = \mu_+T_+ - \mu_-T_- + (F'_+ - F'_-) = \mu_+T_+ - (\mu_+ - \kappa)T_- + (F'_+ - F'_-)$$

$$= \mu_+T_a + \kappa T_- + (F'_+ - F'_-) = \mu_+T_a + \frac{\kappa}{2}(T_s - T_a) + (F'_+ - F'_-)$$

$$= \frac{\kappa}{2}T_s + \frac{\hat{\lambda}}{2}T_a + (F'_+ - F'_-) = \frac{\kappa}{2}T_s + \frac{\hat{\lambda}}{2}T_a + \frac{\hat{\lambda}+2\gamma'_d}{\kappa}F' \therefore$$

$$N = \frac{1}{2}\left(\hat{\lambda}\mathrm{C}_s + \kappa\mathrm{C}_a\right)T_s + \frac{1}{2}\left(\hat{\lambda}\mathrm{C}_a + \kappa\mathrm{C}_s\right)T_a + \left(\mathrm{C}_s + \mathrm{C}_a\frac{\hat{\lambda}+2\gamma'_d}{\kappa}\right)F'$$

27

581    Further expanding the coefficients

$$
\hat{\lambda} C_s + \kappa C_a = \hat{\lambda} C_u - \frac{C_d}{2\varepsilon\gamma'}(\hat{\lambda}^2 + 2\gamma'_d \hat{\lambda} - \kappa^2) = \hat{\lambda} C_u - \frac{C_d}{2\varepsilon\gamma'}(\hat{\lambda}^2 + 2\gamma'_d \hat{\lambda} - \hat{\lambda}^2 - 4\gamma'_d \lambda')
$$

$$
= 2\frac{C_u}{\varepsilon}\left(\lambda' + \frac{\varepsilon - 1}{2}\hat{\lambda}\right)
$$

$$
\hat{\lambda} C_a + \kappa C_s = \kappa C_u - \frac{C_d}{2\varepsilon\gamma'}(\kappa\hat{\lambda} + 2\gamma'_d \kappa - \kappa\hat{\lambda}) = \kappa C_u - \frac{C_u}{\varepsilon}\kappa = \kappa\frac{C_u}{\varepsilon}(\varepsilon - 1)
$$

$$
C_s + C_a\frac{\hat{\lambda} + 2\gamma'_d}{\kappa} = C_u - \frac{C_d}{2\varepsilon\gamma'}(\hat{\lambda} + 2\gamma'_d - \hat{\lambda} - 2\gamma'_d) = C_u
$$

587    then the imbalance is

$$
N = \frac{C_u}{\varepsilon}\left[\varepsilon F' + \left(\lambda' + \frac{\varepsilon - 1}{2}\hat{\lambda}\right)T_s + \kappa\frac{\varepsilon - 1}{2}T_a\right] \tag{A14}
$$

590    From here, I derive the slope of a $NT$−diagram. In such a diagram, $N$ is plotted versus $T_u$. If we

591    naïvely take the partial derivative of equation (A14) with respect to $T_u$, we will arrive to a constant

592    slope. This is contrary to the evidence that it will change with time. An $NT$−diagram is one

593    projection of the phase space of the system. Then, the $NT$−diagram slope does not only depend on

594    how $N$ varies with $T_u$. It is a comparison of how the changes of $T_u$ are expressed in changes of $N$.

595    Then, the slope is the total derivative $dN/dT_u$. By virtue of the chain rule, $dN/dT_u = \dot{N}(dt/dT_u)$.

596    In a neighborhood where $T_u(t)$ is injective, $dt/dT_u = 1/\dot{T}_u$. Therefore, the slope $dN/dT_u$ is the ratio

597    of two total derivatives: $\dot{N}$ and $\dot{T}_u$.

598    We know that $T_u = T_s$, then $\dot{T}_u = \dot{T}_s$. Therefore, the total derivative of the planetary imbalance is

$$
\dot{N} = (\partial_t N) + (\partial_{T_s} N)\dot{T}_s + (\partial_{T_a} N)\dot{T}_a
$$

601    that is a change depending only on time, a second change depending only on changes of $T_s$ and a

602    third depending on changes of $T_a$. Therefore, the ratio of total derivative of planetary imbalance

603    and total derivative of $T_u$ is

$$
\frac{\dot{N}}{\dot{T}_u} = (\partial_t N)\frac{1}{\dot{T}_s} + (\partial_{T_s} N) + (\partial_{T_a} N)\frac{\dot{T}_a}{\dot{T}_s}
$$

28

As one can see in the above expression, the ratio includes the derivative of the imbalance with respect to $T_u$ but is not the only contribution. One contribution comes from the explicit dependence on time of $N$ and how it compares with the dependency of $T_u$. The other contribution comes from the antisymmetric mode and how it changes in relation to the symmetric one. From equation (A14), I can write the precise expression of the slope as a factor of $\lambda$.

I multiply equation (A14) by $\lambda/\lambda$ and reorganise.

$$
\frac{\dot{N}}{\dot{T_u}} = \frac{C_u}{\varepsilon}\left[\varepsilon\frac{\dot{F'}}{\dot{T_s}} + \left(\lambda' + \frac{\varepsilon-1}{2}\hat{\lambda}\right) + \kappa\frac{\varepsilon-1}{2}\frac{\dot{T_a}}{\dot{T_s}}\right]\frac{\lambda}{\lambda}
$$

$$
= \left[\frac{C_u}{\lambda}\frac{\dot{F'}}{\dot{T_s}} + \left(\frac{\lambda'}{\varepsilon\lambda'} + \frac{\varepsilon-1}{2\varepsilon}\frac{\hat{\lambda}}{\lambda'}\right) + \frac{\varepsilon-1}{2\varepsilon}\frac{\kappa}{\lambda'}\frac{\dot{T_a}}{\dot{T_s}}\right]\lambda
$$

then we will expand the terms to separate the terms that vanish when $\varepsilon = 1$

$$
\frac{\dot{N}}{\dot{T_u}} = \left\{\frac{C_u}{\lambda}\frac{\dot{F'}}{\dot{T_s}} + \left[\frac{1}{\varepsilon} + \frac{\varepsilon-1}{2\varepsilon}\left(\frac{\lambda'-\varepsilon\gamma'-\gamma_d'}{\lambda'}\right)\right] + \frac{\varepsilon-1}{2\varepsilon}\frac{\kappa}{\lambda'}\frac{\dot{T_a}}{\dot{T_s}}\right\}\lambda
$$

$$
= \left\{\frac{C_u}{\lambda}\frac{\dot{F'}}{\dot{T_s}} + \left[\frac{2}{2\varepsilon} + \frac{\varepsilon-1}{2\varepsilon}\left(1 - \varepsilon\frac{\gamma}{\lambda} - \frac{C_u}{C_d}\frac{\gamma}{\lambda}\right)\right] + \frac{\varepsilon-1}{2\varepsilon}\frac{C_u\kappa}{\lambda}\frac{\dot{T_a}}{\dot{T_s}}\right\}\lambda
$$

$$
= \left[\frac{C_u}{\lambda}\frac{\dot{F'}}{\dot{T_s}} + \frac{\varepsilon+1}{2\varepsilon} - \frac{\varepsilon-1}{2\varepsilon}\left(\varepsilon + \frac{C_u}{C_d}\right)\frac{\gamma}{\lambda} + \frac{\varepsilon-1}{2\varepsilon}\frac{C_u\kappa}{\lambda}\frac{\dot{T_a}}{\dot{T_s}}\right]\lambda
$$

$$
= \left[\frac{C_u}{\lambda}\frac{\dot{F'}}{\dot{T_s}} + \frac{\varepsilon+1}{2\varepsilon} - \frac{\varepsilon-1}{2\varepsilon}\left(\varepsilon + \frac{C_u}{C_d}\right)\frac{\gamma}{\lambda} + \frac{\varepsilon-1}{2\varepsilon}\frac{C_u\kappa}{\lambda}\frac{\dot{T_a}}{\dot{T_s}}\right]\lambda
$$

$$
= \left\{\frac{C_u}{\lambda}\frac{\dot{F'}}{\dot{T_s}} + \frac{\varepsilon+1}{2\varepsilon} - \frac{\varepsilon-1}{2\varepsilon\lambda}\left[\left(\varepsilon + \frac{C_u}{C_d}\right)\gamma - C_u\kappa\frac{\dot{T_a}}{\dot{T_s}}\right]\right\}\lambda
$$

$$
= \left\{\frac{C_u}{\lambda}\frac{\dot{F'}}{\dot{T_s}} + \frac{\varepsilon+1}{2\varepsilon} - \frac{\varepsilon-1}{2\varepsilon\lambda}C_u\kappa\left[\left(\varepsilon + \frac{C_u}{C_d}\right)\frac{\gamma}{C_u\kappa} - \frac{\dot{T_a}}{\dot{T_s}}\right]\right\}\lambda
$$

$$
= \left\{\frac{C_u}{\lambda}\frac{\dot{F'}}{\dot{T_s}} + \frac{\varepsilon+1}{2\varepsilon} - \frac{\varepsilon-1}{2\varepsilon}\frac{C_u\kappa}{\lambda}\left[\left(\varepsilon + \frac{C_u}{C_d}\right)\frac{\gamma}{C_u\kappa} - \frac{\dot{T_a}}{\dot{T_s}}\right]\right\}\lambda
$$

$$
= \left\{-\frac{C_u}{|\lambda|}\frac{\dot{F'}}{\dot{T_s}} + \frac{\varepsilon+1}{2\varepsilon} + \frac{\varepsilon-1}{2\varepsilon}\frac{C_u\kappa}{|\lambda|}\left[\left(\varepsilon + \frac{C_u}{C_d}\right)\frac{\gamma}{C_u\kappa} - \frac{\dot{T_a}}{\dot{T_s}}\right]\right\}\lambda
$$

$$
\frac{\dot{N}}{\dot{T_u}} = \left\{-\frac{C_u}{|\lambda|}\frac{\dot{F'}}{\dot{T_s}} + \frac{\varepsilon+1}{2\varepsilon}\left(1 + \frac{\varepsilon-1}{\varepsilon+1}\frac{C_u\kappa}{|\lambda|}\left[\left(\varepsilon + \frac{C_u}{C_d}\right)\frac{\gamma}{C_u\kappa} - \frac{\dot{T_a}}{\dot{T_s}}\right]\right)\right\}\lambda \tag{A15}
$$

The term in square brackets in equation (A15) is the key term that provides a $NT-$diagram with evolving slope when the forcing is constant. The second part of this term provides the temporal

29

630 evolution, whereas the first part is a constant term that sets the base enhancement of the slope.

631 Interestingly, this first part contains in particular the thermal capacities of the system.

632 If I rewrite this first part of the square-brackets term, the terms are shown clearly

633
634
$$\frac{\dot{N}}{\dot{T}_{\mathrm{u}}} = \left\{ -\frac{C_{\mathrm{u}}}{|\lambda|}\frac{\dot{F}'}{\dot{T}_s} + \frac{\varepsilon+1}{2\varepsilon} + \frac{\varepsilon-1}{2\varepsilon}\frac{C_{\mathrm{u}}\kappa}{|\lambda|}\left[\left(\frac{\varepsilon}{C_{\mathrm{u}}} + \frac{1}{C_{\mathrm{d}}}\right)\frac{\gamma}{\kappa} - \frac{\dot{T}_a}{\dot{T}_s}\right]\right\}\lambda \tag{A16}$$

635 Now in the first part it is the sum of the inverse of the thermal capacities as if we have an electrical

636 circuit with capacitors in series. Having such a term in the equation for the slope favors the physical

637 interpretation in terms of thermal capacities, instead of variable feedback mechanisms. The time-

638 evolving ratio term in the second part, that represents the dynamics of the atmosphere-ocean

639 coupling, only strengthens this interpretation.

640 As a corollary, if the forcing is constant and $\varepsilon \to 1$, then we recover the classical linear dependence

641 of the imbalance on $T_{\mathrm{u}}$

642
643
$$\lim_{\varepsilon\to 1}\frac{\dot{N}}{\dot{T}_{\mathrm{u}}} = \lambda, \, F = \text{const}$$

644 **Symmetric and antisymmetric modes**

645 From equations (A13), we see that the symmetric and antisymmetric modes are the basis for

646 the description of the solutions. Thus, let us give some explicit expression for the symmetric and

647 antisymmetric modes.

30

648 From equation (A12) and the equations for the initial values and the forcing, I can write more

649 explicitly the solution

$$
T_\pm = \left(T_{\pm,0} + \int_{t_0}^{t} F'_\pm e^{-\mu_\pm(\tau-t_0)}\mathrm{d}\tau\right) e^{\mu_\pm(t-t_0)}
$$

$$
= \left(\pm\frac{1}{\mu_+ - \mu_-}[(\mu_\pm + \gamma'_d)T_{u,0} + \varepsilon\gamma'T_{d,0}] \pm \frac{\mu_\pm + \gamma'_d}{\mu_+ - \mu_-}\int_{t_0}^{t} F' e^{-\mu_\pm(\tau-t_0)}\mathrm{d}\tau\right) e^{\mu_\pm(t-t_0)}
$$

$$
= \pm\frac{e^{(\hat\lambda/2)(t-t_0)}}{\mu_+ - \mu_-}\left[(\mu_\pm + \gamma'_d)T_{u,0} + \varepsilon\gamma'T_{d,0} + (\mu_\pm + \gamma'_d)\int_{t_0}^{t} F' e^{-\mu_\pm(\tau-t_0)}\mathrm{d}\tau\right] e^{\pm(\kappa/2)(t-t_0)}
$$

$$
= \pm\frac{e^{(\hat\lambda/2)(t-t_0)}}{\mu_+ - \mu_-}\left[\frac{\hat\lambda \pm \kappa + 2\gamma'_d}{2}T_{u,0} + \frac{2\varepsilon\gamma'}{2}T_{d,0} + \frac{\hat\lambda \pm \kappa + 2\gamma'_d}{2}\int_{t_0}^{t} F' e^{-\mu_\pm(\tau-t_0)}\mathrm{d}\tau\right] e^{\pm(\kappa/2)(t-t_0)}
$$

$$
= \pm\frac{e^{(\hat\lambda/2)(t-t_0)}}{2(\mu_+ - \mu_-)}\left[(\hat\lambda + 2\gamma'_d)T_{u,0} + 2\varepsilon\gamma'T_{d,0} \pm \kappa T_{u,0} + (\hat\lambda + 2\gamma'_d \pm \kappa)\int_{t_0}^{t} F' e^{-\mu_\pm(\tau-t_0)}\mathrm{d}\tau\right] e^{\pm(\kappa/2)(t-t_0)}
$$

656 Now that I have a more explicit expression, I write the modes

657
$$
T_+ \pm T_- =
$$

$$
\frac{e^{(\hat\lambda/2)(t-t_0)}}{2(\mu_+ - \mu_-)}\left[(\hat\lambda + 2\gamma'_d)T_{u,0} + 2\varepsilon\gamma'T_{d,0} + \kappa T_{u,0} + (\hat\lambda + 2\gamma'_d + \kappa)\int_{t_0}^{t} F' e^{-\mu_+(\tau-t_0)}\mathrm{d}\tau\right] e^{(\kappa/2)(t-t_0)}
$$

$$
\mp\frac{e^{(\hat\lambda/2)(t-t_0)}}{2(\mu_+ - \mu_-)}\left[(\hat\lambda + 2\gamma'_d)T_{u,0} + 2\varepsilon\gamma'T_{d,0} - \kappa T_{u,0} + (\hat\lambda + 2\gamma'_d - \kappa)\int_{t_0}^{t} F' e^{-\mu_-(\tau-t_0)}\mathrm{d}\tau\right] e^{-(\kappa/2)(t-t_0)}
$$

$$
= \frac{e^{(\hat\lambda/2)(t-t_0)}}{\mu_+ - \mu_-}\left\{\left[(\hat\lambda + 2\gamma'_d)T_{u,0} + 2\varepsilon\gamma'T_{d,0}\right]\frac{e^{(\kappa/2)(t-t_0)} \mp e^{-(\kappa/2)(t-t_0)}}{2}\right.
$$

$$
+ \kappa T_{u,0}\frac{e^{(\kappa/2)(t-t_0)} \pm e^{-(\kappa/2)(t-t_0)}}{2}
$$

$$
+ \frac{\hat\lambda + 2\gamma'_d}{2}\left[e^{(\kappa/2)(t-t_0)}\int_{t_0}^{t} F' e^{-\mu_+(\tau-t_0)}\mathrm{d}\tau \mp e^{-(\kappa/2)(t-t_0)}\int_{t_0}^{t} F' e^{-\mu_-(\tau-t_0)}\mathrm{d}\tau\right]
$$

$$
\left. + \frac{\kappa}{2}\left[e^{(\kappa/2)(t-t_0)}\int_{t_0}^{t} F' e^{-\mu_+(\tau-t_0)}\mathrm{d}\tau \pm e^{-(\kappa/2)(t-t_0)}\int_{t_0}^{t} F' e^{-\mu_-(\tau-t_0)}\mathrm{d}\tau\right]\right\}
$$

665 The last two terms inside the curly brackets have a similar form as the combinations of exponential

666 functions in the first two terms. These combinations of exponential functions are hyperbolic

667 functions which can simplify the expressions of the solutions. I would want such a representation

668 but a problem is there: the integrals are not the same, therefore I cannot factorise them together.

669 Notwithstanding, from the definition of hyperbolic sine and cosine functions, I can write $e^{\pm x} =$

31

670 $\cosh x \pm \sinh x$. The factors within square brackets in the last two terms can be thought as $e^x I_+ \pm$

671 $e^{-x} I_-$, where $I_\pm$ are the corresponding integrals. Using the expression of the exponential function in

672 terms of the hyperbolic functions, I expand $e^x I_+ \pm e^{-x} I_- = (\cosh x + \sinh x) I_+ \pm (\cosh x - \sinh x) I_- =$

673 $(I_+ \pm I_-) \cosh x + (I_+ \mp I_-) \sinh x$. Then, I overcome the limitation and now the two terms are written

674 with hyperbolic functions. The coefficients of the hyperbolic functions are simple combinations

675 of the integrals which can be also expanded easily. I do that now

676
$$I_+ + I_- = \int_{t_0}^t F' e^{-\mu_+ (\tau - t_0)} d\tau + \int_{t_0}^t F' e^{-\mu_- (\tau - t_0)} d\tau = \int_{t_0}^t F'[e^{-\mu_+ (\tau - t_0)} + e^{-\mu_- (\tau - t_0)}] d\tau$$

677
$$= \int_{t_0}^t F' e^{-(\hat{\lambda}/2)(\tau - t_0)} [e^{-(\kappa/2)(\tau - t_0)} + e^{(\kappa/2)(\tau - t_0)}] d\tau$$

678
$$= 2 \int_{t_0}^t F' e^{-(\hat{\lambda}/2)(\tau - t_0)} \cosh \left[ \frac{\kappa}{2} (\tau - t_0) \right] d\tau$$

679
$$I_+ - I_- = \int_{t_0}^t F' e^{-\mu_+ (\tau - t_0)} d\tau - \int_{t_0}^t F' e^{-\mu_- (\tau - t_0)} d\tau = \int_{t_0}^t F'[e^{-\mu_+ (\tau - t_0)} - e^{-\mu_- (\tau - t_0)}] d\tau$$

680
$$= \int_{t_0}^t F' e^{-(\hat{\lambda}/2)(\tau - t_0)} [e^{-(\kappa/2)(\tau - t_0)} - e^{(\kappa/2)(\tau - t_0)}] d\tau$$

681
$$= -2 \int_{t_0}^t F' e^{-(\hat{\lambda}/2)(\tau - t_0)} \sinh \left[ \frac{\kappa}{2} (\tau - t_0) \right] d\tau$$

682

683 If one collects terms corresponding to each hyperbolic function in the former expressions for the

684 normal modes, obtains the following

685
$$T_s = \frac{e^{(\hat{\lambda}/2)(t - t_0)}}{\kappa} \left\{ C_1 \cosh \left[ \frac{\kappa}{2} (t - t_0) \right] + C_2 \sinh \left[ \frac{\kappa}{2} (t - t_0) \right] \right\} \tag{A17}$$

686
$$T_a = \frac{e^{(\hat{\lambda}/2)(t - t_0)}}{\kappa} \left\{ C_2 \cosh \left[ \frac{\kappa}{2} (t - t_0) \right] + C_1 \sinh \left[ \frac{\kappa}{2} (t - t_0) \right] \right\} \tag{A18}$$

687

32

where

$$C_1 = \kappa T_{u,0}$$

$$- (\hat{\lambda} + 2\gamma_d') \int_{t_0}^{t} F' e^{-(\hat{\lambda}/2)(\tau - t_0)} \sinh\left[\frac{\kappa}{2}(\tau - t_0)\right] d\tau + \kappa \int_{t_0}^{t} F' e^{-(\hat{\lambda}/2)(\tau - t_0)} \cosh\left[\frac{\kappa}{2}(\tau - t_0)\right] d\tau$$

$$C_2 = (\hat{\lambda} + 2\gamma_d')T_{u,0} + 2\varepsilon\gamma_d' T_{d,0}$$

$$+ (\hat{\lambda} + 2\gamma_d') \int_{t_0}^{t} F' e^{-(\hat{\lambda}/2)(\tau - t_0)} \cosh\left[\frac{\kappa}{2}(\tau - t_0)\right] d\tau - \kappa \int_{t_0}^{t} F' e^{-(\hat{\lambda}/2)(\tau - t_0)} \sinh\left[\frac{\kappa}{2}(\tau - t_0)\right] d\tau$$

These expressions for the normal modes are quite elegant, and the coefficients $C_i$ summarize all the information from the initial conditions and the forcing. The initial condition terms in the $C_i$ correspond to the non-forced response of the system, while the part that is forcing-dependent corresponds to the forced response of the system.

**Forced response to constant forcing**

If $F' = F_c' \neq 0$ for $t > t_0$ with $F_c'$ constant and $T_{u,0}, T_{d,0} = 0$ for $t = t_0$, then

$$C_1 = F_c' \left\{ -(\hat{\lambda} + 2\gamma_d') \int_{t_0}^{t} e^{-(\hat{\lambda}/2)(\tau - t_0)} \sinh\left[\frac{\kappa}{2}(\tau - t_0)\right] d\tau + \kappa \int_{t_0}^{t} e^{-(\hat{\lambda}/2)(\tau - t_0)} \cosh\left[\frac{\kappa}{2}(\tau - t_0)\right] d\tau \right\}$$

$$C_2 = F_c' \left\{ (\hat{\lambda} + 2\gamma_d') \int_{t_0}^{t} e^{-(\hat{\lambda}/2)(\tau - t_0)} \cosh\left[\frac{\kappa}{2}(\tau - t_0)\right] d\tau - \kappa \int_{t_0}^{t} e^{-(\hat{\lambda}/2)(\tau - t_0)} \sinh\left[\frac{\kappa}{2}(\tau - t_0)\right] d\tau \right\}$$

where the integrals are easily computed

$$\int_{t_0}^{t} e^{-(\hat{\lambda}/2)(\tau - t_0)} \sinh\left[\frac{\kappa}{2}(\tau - t_0)\right] d\tau = \frac{e^{-(\hat{\lambda}/2)(t - t_0)}}{\lambda'\gamma_d'} \left\{ \frac{\kappa}{2} \cosh\left[\frac{\kappa}{2}(t - t_0)\right] + \frac{\hat{\lambda}}{2} \sinh\left[\frac{\kappa}{2}(t - t_0)\right] \right\} - \frac{\kappa}{2\lambda'\gamma_d'}$$

$$\int_{t_0}^{t} e^{-(\hat{\lambda}/2)(\tau - t_0)} \cosh\left[\frac{\kappa}{2}(\tau - t_0)\right] d\tau = \frac{e^{-(\hat{\lambda}/2)(t - t_0)}}{\lambda'\gamma_d'} \left\{ \frac{\hat{\lambda}}{2} \cosh\left[\frac{\kappa}{2}(t - t_0)\right] + \frac{\kappa}{2} \sinh\left[\frac{\kappa}{2}(t - t_0)\right] \right\} - \frac{\hat{\lambda}}{2\lambda'\gamma_d'}$$

and, upon reduction, the $C_i$ are

$$C_1 = \frac{F_c'}{\lambda'} e^{-(\hat{\lambda}/2)(\tau - t_0)} \left\{ -\kappa \cosh\left[\frac{\kappa}{2}(t - t_0)\right] + (2\lambda' - \hat{\lambda}) \sinh\left[\frac{\kappa}{2}(t - t_0)\right] + \kappa e^{(\hat{\lambda}/2)(t - t_0)} \right\}$$

$$C_2 = \frac{F_c'}{\lambda'} e^{-(\hat{\lambda}/2)(\tau - t_0)} \left\{ -(2\lambda' - \hat{\lambda}) \cosh\left[\frac{\kappa}{2}(t - t_0)\right] + \kappa \sinh\left[\frac{\kappa}{2}(t - t_0)\right] + (2\lambda' - \hat{\lambda}) e^{(\hat{\lambda}/2)(t - t_0)} \right\}$$

33

711 with these expressions is easy to evaluate the terms inside the curly brackets in equations (A17)

712 and (A18) and the symmetric and antisymmetric modes are (for $t \geq t_0$)

$$T_s = \frac{F_c}{\lambda} \left\{ e^{(\hat{\lambda}/2)(t-t_0)} \left( \cosh\left[\frac{\kappa}{2}(t-t_0)\right] + \frac{2\lambda' - \hat{\lambda}}{\kappa} \sinh\left[\frac{\kappa}{2}(t-t_0)\right] \right) - 1 \right\} \tag{A19}$$

$$T_a = \frac{F_c}{\lambda} \left\{ e^{(\hat{\lambda}/2)(t-t_0)} \left( \frac{2\lambda' - \hat{\lambda}}{\kappa} \cosh\left[\frac{\kappa}{2}(t-t_0)\right] + \sinh\left[\frac{\kappa}{2}(t-t_0)\right] \right) - \frac{2\lambda' - \hat{\lambda}}{\kappa} \right\} \tag{A20}$$

716 where $F_c' := F_c/C_u$. I can also obtain the explicit time derivatives of both modes. We take the time

717 derivative both equations (A19) and (A20)

$$
\begin{aligned}
\dot{T}_s &= \frac{F_c}{\lambda} e^{(\hat{\lambda}/2)(t-t_0)} \left\{ \frac{\hat{\lambda}}{2} \left( \cosh\left[\frac{\kappa}{2}(t-t_0)\right] + \frac{2\lambda' - \hat{\lambda}}{\kappa} \sinh\left[\frac{\kappa}{2}(t-t_0)\right] \right) \right.\\
&\qquad\qquad\qquad \left. + \frac{\kappa}{2} \left( \frac{2\lambda' - \hat{\lambda}}{\kappa} \cosh\left[\frac{\kappa}{2}(t-t_0)\right] + \sinh\left[\frac{\kappa}{2}(t-t_0)\right] \right) \right\}\\
&= \frac{F_c}{\lambda} e^{(\hat{\lambda}/2)(t-t_0)} \left\{ \lambda' \cosh\left[\frac{\kappa}{2}(t-t_0)\right] + \frac{\lambda'\hat{\lambda} + 2\gamma_d'\lambda'}{\kappa} \sinh\left[\frac{\kappa}{2}(t-t_0)\right] \right\}\\
&= \frac{F_c}{C_u} e^{(\hat{\lambda}/2)(t-t_0)} \left\{ \cosh\left[\frac{\kappa}{2}(t-t_0)\right] + \frac{\hat{\lambda} + 2\gamma_d'}{\kappa} \sinh\left[\frac{\kappa}{2}(t-t_0)\right] \right\}
\end{aligned}
$$

$$
\begin{aligned}
\dot{T}_a &= \frac{F_c}{\lambda} e^{(\hat{\lambda}/2)(t-t_0)} \left\{ \frac{\hat{\lambda}}{2} \left( \frac{2\lambda' - \hat{\lambda}}{\kappa} \cosh\left[\frac{\kappa}{2}(t-t_0)\right] + \sinh\left[\frac{\kappa}{2}(t-t_0)\right] \right) \right.\\
&\qquad\qquad\qquad \left. + \frac{\kappa}{2} \left( \cosh\left[\frac{\kappa}{2}(t-t_0)\right] + \frac{2\lambda' - \hat{\lambda}}{\kappa} \sinh\left[\frac{\kappa}{2}(t-t_0)\right] \right) \right\}\\
&= \frac{F_c}{\lambda} e^{(\hat{\lambda}/2)(t-t_0)} \left\{ \frac{\lambda'\hat{\lambda} + 2\gamma_d'\lambda'}{\kappa} \cosh\left[\frac{\kappa}{2}(t-t_0)\right] + \lambda' \sinh\left[\frac{\kappa}{2}(t-t_0)\right] \right\}\\
&= \frac{F_c}{C_u} e^{(\hat{\lambda}/2)(t-t_0)} \left\{ \frac{\hat{\lambda} + 2\gamma_d'}{\kappa} \cosh\left[\frac{\kappa}{2}(t-t_0)\right] + \sinh\left[\frac{\kappa}{2}(t-t_0)\right] \right\}
\end{aligned}
$$

727 I present both results jointly to show the simplicity of the derivatives

$$\dot{T}_s = \frac{F_c}{C_u} e^{(\hat{\lambda}/2)(t-t_0)} \left\{ \cosh\left[\frac{\kappa}{2}(t-t_0)\right] + \frac{\hat{\lambda} + 2\gamma_d'}{\kappa} \sinh\left[\frac{\kappa}{2}(t-t_0)\right] \right\}$$

$$\dot{T}_a = \frac{F_c}{C_u} e^{(\hat{\lambda}/2)(t-t_0)} \left\{ \frac{\hat{\lambda} + 2\gamma_d'}{\kappa} \cosh\left[\frac{\kappa}{2}(t-t_0)\right] + \sinh\left[\frac{\kappa}{2}(t-t_0)\right] \right\}$$

34

With these derivatives, I can calculate the ratio of the antisymmetric mode derivative to the symmetric one that appears in equation (A15)

$$\frac{\dot{T}_a}{\dot{T}_s} = \frac{\frac{\hat{\lambda}+2\gamma'_d}{\kappa}\cosh\left[\frac{\kappa}{2}(t-t_0)\right] + \sinh\left[\frac{\kappa}{2}(t-t_0)\right]}{\cosh\left[\frac{\kappa}{2}(t-t_0)\right] + \frac{\hat{\lambda}+2\gamma'_d}{\kappa}\sinh\left[\frac{\kappa}{2}(t-t_0)\right]}$$

$$= \frac{\frac{\hat{\lambda}+2\gamma'_d}{\kappa} + \tanh\left[\frac{\kappa}{2}(t-t_0)\right]}{1 + \frac{\hat{\lambda}+2\gamma'_d}{\kappa}\tanh\left[\frac{\kappa}{2}(t-t_0)\right]}$$

Formally, above result have the alternative form

$$\frac{\dot{T}_a}{\dot{T}_s} = \tanh\left[\frac{\kappa}{2}(t-t_0) + \operatorname{arctanh}\left(\frac{\hat{\lambda}+2\gamma'_d}{\kappa}\right)\right]$$

This is possible only if $\left|(\hat{\lambda}+2\gamma'_d)/\kappa\right| \leq 1$. Let us prove that in our case this follows

$$\left|\frac{\hat{\lambda}+2\gamma'_d}{\kappa}\right| \leq 1$$

$$\frac{\hat{\lambda}^2 + 4\gamma'_d\hat{\lambda} + 4\gamma'^2_d}{\hat{\lambda}^2 + 4\gamma'_d\lambda'} \leq 1$$

$$\hat{\lambda}^2 + 4\gamma'_d\hat{\lambda} + 4\gamma'^2_d \leq \hat{\lambda}^2 + 4\gamma'_d\lambda'$$

$$\hat{\lambda} + \gamma'_d \leq \lambda'$$

$$-\varepsilon\gamma' \leq 0$$

the last inequality is always true, since $\varepsilon, \gamma'$ are positive constants. Thus,

$$\frac{\dot{T}_a}{\dot{T}_s} = \tanh\left[\frac{\kappa}{2}(t-t_0) + \operatorname{arctanh}\left(\frac{\hat{\lambda}+2\gamma'_d}{\kappa}\right)\right] \tag{A21}$$

Equation (A21) is an hyperbolic tangent that grows from -1 to 1 in a sigmoidal fashion. It has a scaling factor that determines how fast it goes from -1 to 1. It also has a shift that sets where the hyperbolic tangent will cross zero. Both the scaling and shift depend on the thermal and radiative parameters of the system. Since the shift is negative, after the initial forcing the deep ocean (that depends on the antisymmetric mode) warms up slower than the upper ocean. At a latter time, the

35

754    ratio becomes positive and the contrary happens. The time at which the sign reverses is

$$t_1 = t_0 + \frac{2}{\kappa} \operatorname{arctanh} \left| \frac{\hat{\lambda} + 2\gamma'_d}{\kappa} \right|$$

## Variation of the climate feedback parameter

758    With the solution shown before, the $NT-$diagram has a slope

$$\frac{\dot{N}}{\dot{T}_u} = \frac{\varepsilon+1}{2\varepsilon} \left( 1 + \frac{\varepsilon-1}{\varepsilon+1} \frac{C_u \kappa}{|\lambda|} \left[ \left( \varepsilon + \frac{C_u}{C_d} \right) \frac{\gamma}{C_u \kappa} - \tanh\left( \frac{\kappa}{2}(t-t_0) + \operatorname{arctanh}\left( \frac{\hat{\lambda}+2\gamma'_d}{\kappa} \right) \right) \right] \right) \lambda \quad \text{(A22)}$$

761    The factor is composed of terms that are positive except for the ratio term coming from equation

762    (A21). The negative ratio for $t \in [t_0, t_1)$ clearly generates a more negative slope, whereas for

763    $t \in (t_1, \infty)$ makes it less negative. At the start one can get the slope

$$\frac{\dot{N}}{\dot{T}_u} = \left( 1 + (\varepsilon-1)\frac{\gamma}{|\lambda|} \right) \lambda, t = t_0$$

766    and at the time of sign reversal

$$\frac{\dot{N}}{\dot{T}_u} = \frac{\varepsilon+1}{2\varepsilon} \left( 1 + \frac{\varepsilon-1}{\varepsilon+1} \left( \varepsilon + \frac{C_u}{C_d} \right) \frac{\gamma}{|\lambda|} \right) \lambda, t = t_1$$

769    After the sign reversal the factor of $\lambda$ will only decrease up to

$$\lim_{t \to \infty} \frac{\dot{N}}{\dot{T}_u} = \frac{\varepsilon+1}{2\varepsilon} \left( 1 + \frac{\varepsilon-1}{\varepsilon+1} \frac{C_u \kappa}{|\lambda|} \left[ \left( \varepsilon + \frac{C_u}{C_d} \right) \frac{\gamma}{C_u \kappa} - 1 \right] \right) \lambda$$

772    Equation (A22) shows the importance of the ratio of the symmetric and antisymmetric modes. Its

773    physical meaning, the relationship between the upper- and deep-ocean warming, sets the strength

774    of the variation of the climate feedback, whereas the constant term sets a base enhancement around

775    which the feedback evolves. The thermal capacities of the system determine this constant term.

## APPENDIX B

## Feedbacks and pattern effect in a non-linear planetary budget

I start with a planetary imbalance considering a variation of the planetary thermal capacity

$$N = (1 - \alpha)S + G - \epsilon\sigma(fT_u)^4 - \dot{C}T_u \tag{B1}$$

where $S$ is the incoming solar short-wave flux at the TOA, $\alpha$ is the planetary albedo, $G$ are the remaining natural and anthropogenic energy fluxes, and the last two terms are the planetary long-wave response and the contribution to the radiative response of a varying thermal capacity. As said in the main text, the ocean circulation and the atmosphere-ocean coupling provide the dynamical component of the thermal capacity.

If I compute the total derivative of $N$ then

$$\dot{N} = \left[(1 - \alpha)\dot{S} + \dot{G}\right] - S\dot{\alpha} - \sigma(fT_u)^4\dot{\epsilon} - 4\epsilon\sigma(fT_u)^3(\dot{f}T_u + f\dot{T}_u) - \dot{C}\dot{T}_u - T_u\ddot{C}$$

$$= \left[(1 - \alpha)\dot{S} + \dot{G}\right] - \mathcal{R}$$

Here we can see the first term is the change from a time-evolving forcing. The rest of the terms, $\mathcal{R}$, are atmospheric feedbacks or the effects of ocean circulation and ocean-atmosphere interaction. The fourth term contains the Planck feedback. Let us compare all the terms of $\mathcal{R}$ in comparison to the Planck feedback term $4\epsilon f\sigma(fT_u)^3\dot{T}_u$

$$\mathcal{R} = S\dot{\alpha} + \sigma(fT_u)^4\dot{\epsilon} + 4\epsilon\sigma(fT_u)^3(\dot{f}T_u + f\dot{T}_u) + \dot{C}\dot{T}_u + T_u\ddot{C}$$

$$= 4\epsilon f\sigma(fT_u)^3\dot{T}_u \left[\frac{S}{4\epsilon f\sigma(fT_u)^3}\frac{\dot{\alpha}}{\dot{T}_u} + \frac{T_u}{4\epsilon}\frac{\dot{\epsilon}}{\dot{T}_u} + \frac{T_u}{f}\frac{\dot{f}}{\dot{T}_u} + 1 + \frac{\dot{C}}{4\epsilon f\sigma(fT_u)^3} + \frac{T_u}{4\epsilon f\sigma(fT_u)^3}\frac{\ddot{C}}{\dot{T}_u}\right]$$

By inserting former expression of $\mathcal{R}$ in the total derivative of the planetary imbalance, reordering and dividing by $\dot{T}_u$, we get the analogous expression for the slope of the $NT-$diagrams

$$\frac{\dot{N}}{\dot{T}_u} = \left[(1 - \alpha)\frac{\dot{S}}{\dot{T}_u} + \frac{\dot{G}}{\dot{T}_u}\right]$$

$$- \left[1 + \frac{S}{4\epsilon f\sigma(fT_u)^3}\frac{\dot{\alpha}}{\dot{T}_u} + \frac{T_u}{4\epsilon}\frac{\dot{\epsilon}}{\dot{T}_u} + \frac{T_u}{f}\frac{\dot{f}}{\dot{T}_u} + \frac{\dot{C}}{4\epsilon f\sigma(fT_u)^3} + \frac{T_u}{4\epsilon f\sigma(fT_u)^3}\frac{\ddot{C}}{\dot{T}_u}\right]4\epsilon f\sigma(fT_u)^3$$

The first contribution in the $\mathcal{R}/\dot{T}_u$ term is 1, representing the Planck feedback. The second contribution is the planetary albedo feedback. It includes the surface albedo feedback as well as

804 the short-wave cloud feedback. The third contribution is the emissivity feedback, to which mainly

805 contributes the traditional water-vapor feedback. The fourth contribution is a representation of the

806 lapse-rate feedback. The fifth and sixth contributions are not atmospheric feedbacks but the effect

807 of the evolving planetary thermal capacity provided by the atmosphere-ocean interaction and the

808 ocean circulation.

809     Both the fifth and sixth contributions measure the effect of a changing planetary thermal capacity.

810 The fifth term should be positive but reduces its contribution towards the equilibrium in view of

811 the modified two-layer model results. In the same context, the sixth contribution should change

812 sign, in analogy to the linearized model results.

## References

814 Andrews, T., J. M. Gregory, M. J. Webb, and K. E. Taylor, 2012: Forcing, feedbacks and climate

815     sensitivity in CMIP5 coupled atmosphere-ocean climate models. *Geophys. Res. Lett.*, **39 (9)**,

816     L09 712, doi:10.1029/2012GL051607.

817 Armour, K. C., C. M. Bitz, and G. H. Roe, 2013: Time-Varying Climate Sensitivity from Regional

818     Feedbacks. *J. Climate*, **26 (13)**, 4518–4534, doi:10.1175/JCLI-D-12-00544.1.

819 Ceppi, P., and J. M. Gregory, 2017: Relationship of tropospheric stability to climate sensitivity

820     and Earth's observed radiation budget. *Proc. Natl. Acad. Sci. (USA)*, **114 (50)**, 13 126–13 131,

821     doi:10.1073/pnas.1714308114.

822 Geoffroy, O., D. Saint-Martin, G. Bellon, A. Voldoire, D. J. L. Olivié, and S. Tytéca, 2013a:

823     Transient Climate Response in a Two-Layer Energy-Balance Model. Part II: Representation of

824     the Efficacy of Deep-Ocean Heat Uptake and Validation for CMIP5 AOGCMs. *J. Climate*, **26 (6)**,

825     1859–1876, doi:10.1175/JCLI-D-12-00196.1.

826 Geoffroy, O., D. Saint-Martin, D. J. L. Olivié, A. Voldoire, G. Bellon, and S. Tytéca, 2013b:

827 Transient Climate Response in a Two-Layer Energy-Balance Model. Part I: Analytical Solution

828 and Parameter Calibration Using CMIP5 AOGCM Experiments. *J. Climate*, **26 (6)**, 1841–1857,

829 doi:10.1175/JCLI-D-12-00195.1.

830 Gregory, J. M., R. J. Stouffer, S. C. B. Raper, P. A. Stott, and N. A. Rayner, 2002: An Ob-

831 servationally Based Estimate of the Climate Sensitivity. *J. Climate*, **15 (22)**, 3117–3121, doi:

832 10.1175/1520-0442(2002)015<3117:AOBEOT>2.0.CO;2.

833 Grose, M. R., J. Gregory, R. Colman, and T. Andrews, 2018: What Climate Sensitiv-

834 ity Index Is Most Useful for Projections? *Geophys. Res. Lett.*, **45 (3)**, 1559–1566, doi:

835 10.1002/2017GL075742.

836 Hawkins, E., and R. Sutton, 2009: The Potential to Narrow Uncertainty in Regional Climate

837 Predictions. *Bull. Amer. Meteor. Soc.*, **90 (8)**, 1095–1107, doi:10.1175/2009BAMS2607.1.

838 Jiménez-de-la-Cuesta, D., and T. Mauritsen, 2019: Emergent constraints on Earth's transient and

839 equilibrium response to doubled $CO_2$ from post-1970s warming. *Nat. Geosci.*, **12 (11)**, 902–905,

840 doi:10.1038/s41561-019-0463-y.

841 Kiehl, J., 2007: Twentieth century climate model response and climate sensitivity. *Geophys. Res.*

842 *Lett.*, **34 (22)**, L22 710, doi:10.1029/2007GL031383.

843 Mauritsen, T., 2016: Clouds cooled the Earth. *Nat. Geosci.*, **9 (12)**, 865–867, doi:10.1038/

844 ngeo2838.

845 Meraner, K., T. Mauritsen, and A. Voigt, 2013: Robust increase in equilibrium climate sensitivity

846 under global warming. *Geophys. Res. Lett.*, **40 (2)**, 5944–5948, doi:10.1002/2013GL058118.

847  Rohrschneider, T., B. Stevens, and T. Mauritsen, 2019: On simple representations of the climate response to external radiative forcing. *Climate Dyn.*, **3 (5-6)**, 3131–3145, doi: 10.1007/s00382-019-04686-4.

850  Senior, C. A., and J. F. B. Mitchell, 2000: The time-dependence of climate sensitivity. *Geophys. Res. Lett.*, **27 (17)**, 2685–2688, doi:10.1029/2000GL011373.

852  Stevens, B., S. C. Sherwood, S. Bony, and M. J. Webb, 2016: Prospects for narrowing bounds on Earth's equilibrium climate sensitivity. *Earths Future*, **4 (11)**, 512–522, doi: 10.1002/2016EF000376.

855  Winton, M., K. Takahashi, and I. M. Held, 2010: Importance of Ocean Heat Uptake Efficacy to Transient Climate Change. *J. Climate*, **23 (9)**, 2333–2344, doi:10.1175/2009JCLI3139.1.

857  Zhou, C., M. D. Zelinka, and S. A. Klein, 2016: Impact of decadal cloud variations on the Earth's energy budget. *Nat. Geosci.*, **9 (12)**, 871–874, doi:10.1038/ngeo2828.
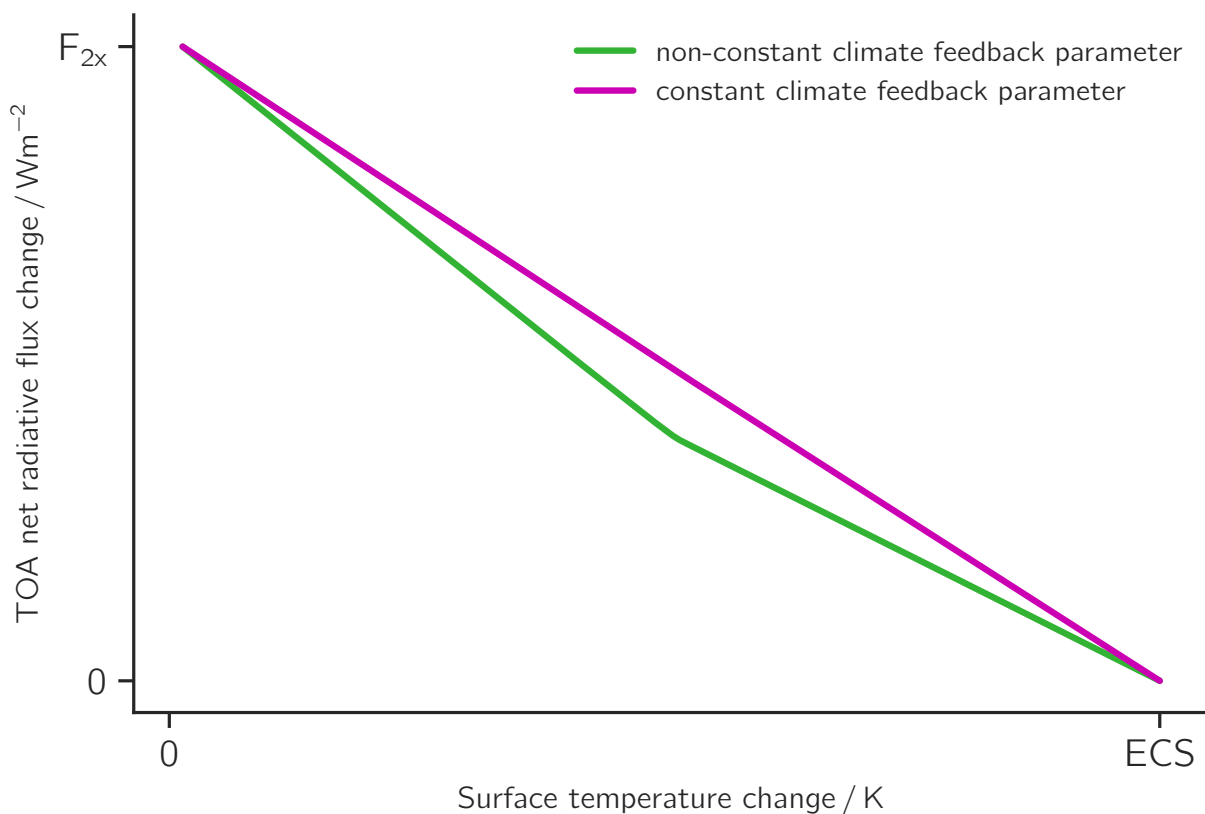
## LIST OF FIGURES

877    FIG. 1. Schematic representation of an $NT-$diagram for constant forcing due to a doubling of the atmospheric

878    carbon dioxide concentration ($F_{2x}$). Magenta line represents the relationship between the TOA net radiatiave

879    flux change with the surface temperature change if the feedback mechanisms on surface warming were constant

880    (constant slope). Green line shows the case found in most models, where the slope varies throughout the process.

881    Given that most models are not run until equilibrium, the evolving slope introduces considerable uncertainty in

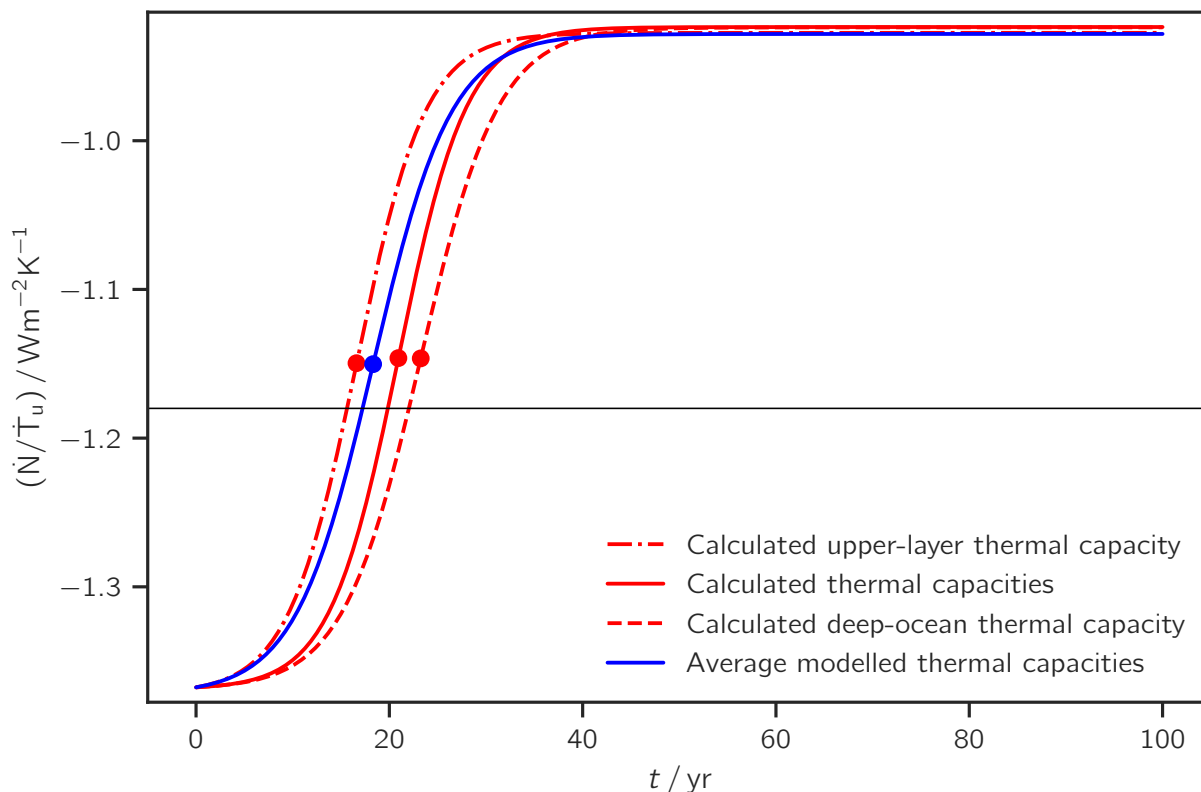882    the equilibrium climate sensitivity (ECS) estimates

FIG. 2. Evolution of the slope of an $NT-$diagram. Blue solid line, with the average parameters from CMIP5 models obtained by Geoffroy et al. (2013a). Red solid line, with the thermal capacities as calculated by Jiménez-de-la-Cuesta and Mauritsen (2019). Red dashed line, with $C_d$ as in Jiménez-de-la-Cuesta and Mauritsen (2019). Red dash-dotted line, with $C_u$ as in Jiménez-de-la-Cuesta and Mauritsen (2019). Dots represent the slope values when the ratio term $\dot{T}_a/\dot{T}_s$ has the sign reversal. Thin black line is the constant $\lambda = -1.18\,\mathrm{W\,m^{-2}\,K^{-1}}$.
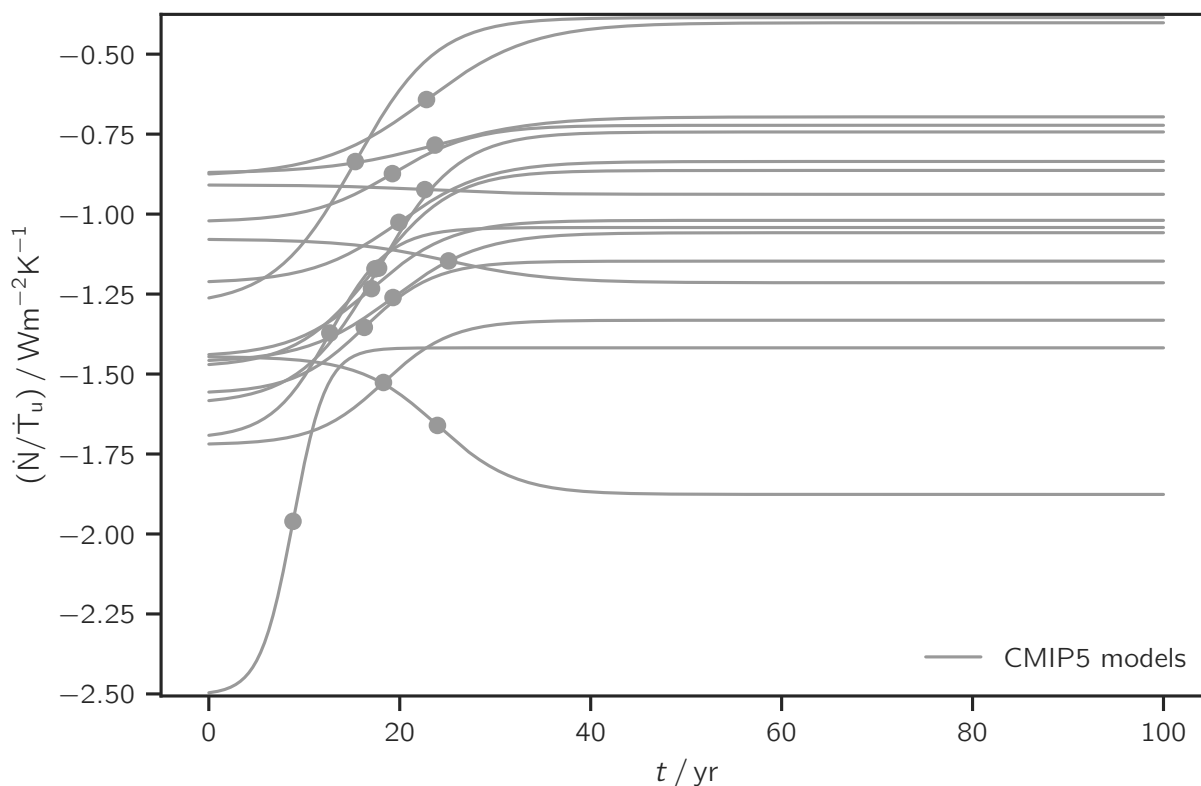
43

888    FIG. 3. Evolution of the slope of the $NT$−diagram. CMIP5 model behaviour using the fitted parameters
889    presented by Geoffroy et al. (2013a). Dots indicate the time of the sign reversal. Note that three models (CNRM-
890    CM5.1, BNU-ESM and INM-CM4) show a steepening slope instead of flattening. For these models, the fitted $\varepsilon$
891    is lesser than one.

44