

Journal of Climate

A Theory for Self-sustained Multicentennial Oscillation of the Atlantic Meridional Overturning Circulation. Part II: Role of Temperature --Manuscript Draft--

Manuscript Number:	JCLI-D-22-0755
Full Title:	A Theory for Self-sustained Multicentennial Oscillation of the Atlantic Meridional Overturning Circulation. Part II: Role of Temperature
Article Type:	Article
Corresponding Author:	Haijun Yang Fudan University Shanghai, CHINA
Corresponding Author's Institution:	Fudan University
First Author:	Kunpeng Yang
Order of Authors:	Kunpeng Yang Haijun Yang Yang Li
Abstract:	<p>In the first part of our research on self-sustained multicentennial oscillation of the Atlantic meridional overturning circulation (AMOC), we used a hemispheric box model considering only the salinity equations. In this follow-up paper, we consider both thermal and saline processes in the box model, so as to investigate the role of temperature in multicentennial AMOC oscillation. The thermal processes exert mainly three effects: shortening the oscillation period, stabilizing subpolar stratification and thus the oscillation system. These three effects are caused by the fast surface temperature restoring process, the stabilizing subpolar temperature stratification, and the negative temperature advection feedback, respectively. Nonlinear restraining effect from enhanced subpolar mixing, or a nonlinear relation between AMOC anomaly and meridional difference of density anomaly, is still needed to realize a self-sustained oscillation, whose mechanism can be generalized as follows: a combination of a linearly growing oscillation dominated by linear advection and a nonlinear restraining process. This study advances the theory reported in the first part of this research. Linear stability analyses reveal that the eigenmode of the system is sensitive to model geometry, flow properties, and meridional differences of sea-surface temperature (SST) and sea-surface salinity (SSS). Our theoretical results suggest that, a smaller (larger) meridional SST (SSS) difference weakens (strengthens) the negative temperature (positive salinity) advection feedback which may lead to a less stable AMOC. Such heuristic findings may be expected in the future due to more intense warming and freshwater hosing at the high latitudes of the Northern Hemisphere.</p>

1 **Replies to Reviewer #1:**

2

3 Thank you very much for all of your constructive comments. We have carefully revised our
4 manuscript based on the advice by you and other reviewers. The following are our point-by-point
5 replies.

6 *This manuscript is a follow up of a previous study published in J. Climate. Both studies used an*
7 *idealized box-model to understand the stability and persistence of a centennial oscillation of the*
8 *AMOC. Whereas the first study focuses solely on the active role of salinity, the present work studies*
9 *the effect of including the temperature as an active variable.*

10 *This work tackles an interesting topic and is a nice and needed follow up of the authors'*
11 *previous work. However, I feel that, unlike the Part I, the theoretical results are not aligned with*
12 *the numerical ones. Hence, I feel that the inconsistency between theory and numerical simulations*
13 *need to be resolved. Also, the fact that the theory developed in Part I is no longer correct when*
14 *temperature is included needs to be discussed.*

15 *Hence, I recommend this work for major revision.*

16 **Responses:** Thank you very much for your invaluable suggestions, which help us improve the
17 manuscript tremendously. Combining the comments from all the reviewers, we have revised the
18 manuscript primarily in these following aspects:

- 19 1) We have completely rewritten the introduction. Coupled modelling studies on multicentennial
20 AMOC oscillation are synthetically reviewed, the inconsistency among their mechanisms and
21 the necessity for theoretical studies are disclosed. Inadequacy of previous theoretical models in
22 accounting for sustainable multicentennial AMOC oscillation is also discussed. Finally, the
23 potential impacts of thermal processes on AMOC oscillation are raised, justifying the inclusion
24 of temperature effects in this study.
- 25 2) In section 2, the choices of parameters are discussed in more detail.
- 26 3) In section 3, we categorize the thermal effects more precisely. We propose that there are mainly
27 three effects when including the temperature equations: (1) increase of the oscillation frequency,
28 (2) stabilization of the overall system, and (3) stabilization of the subpolar stratification. These
29 three effects are attributed to the following three processes, respectively: (1) fast surface
30 temperature restoring, (2) negative temperature advection feedback, and (3) stabilizing subpolar
31 temperature stratification. Now, it is easier to understand that the behaviors of the temperature-
32 salinity system are different from the salinity-only model in LY22.

33 4) In section 4, we more clearly describe the self-sustained oscillation mechanism. In LY22, we
34 proposed that the nonlinear subpolar vertical mixing is crucial for self-sustained oscillation. In
35 the revised manuscript, we further propose that assuming a nonlinear relation between AMOC
36 anomaly and meridional difference of density anomaly can also lead to self-sustained
37 oscillation. We further show that the self-sustained oscillation mechanism not only agrees with
38 that of LY22, but also advances the theory of LY22.

39 5) Most figures are re-plotted.

40

41 **Major Comments:**

42 1. *Inconsistency between theory and numerical simulations*

43 *Overall, the results from analytical approach presented here opposed the one of LY22. In LY22*
44 *the regime parameter is such 4S exhibits an unstable oscillation and 3S a stable one. Hence the*
45 *self-sustained oscillation of [4S + diffusion between box 2 and 3] is interpreted as a mix of the*
46 *unstable oscillation of 4S and the stable one of 3S.*

47 **Responses:** Thank you very much for these comments. Sections 3 and 4 are revised carefully. Our
48 analytical results are not at odds with that of LY22. A more detailed description of the self-
49 sustained oscillation mechanism is provided in section 4, reflecting that the self-sustained
50 oscillation mechanism in the temperature-salinity system agrees well with and even advances that
51 of LY22. Our analytical results and numerical results (of self-sustained oscillation) also match well
52 with each other.

53 Here, we would like to emphasize that, first of all, *in both LY22 and this manuscript, the self-*
54 *sustained oscillation can be never simply interpreted as a mix of an unstable oscillation of 4S (4TS)*
55 *and a stable oscillation of 3S (3TS).* A self-sustained oscillation can never occur in the linear 3S,
56 4S, 3TS, and 4TS systems. The linear stability analyses (no matter using theoretical method or
57 numerical approach) show clearly that in the linear system, there are three kinds of oscillations: the
58 growing oscillation with a positive real part of the eigenvalue, the neutral oscillation with zero real
59 part of the eigenvalue, and the decaying oscillation with a negative real part of the eigenvalue. None
60 of them is a self-sustained oscillation. In different linear systems, the critical value (λ_c) of the linear
61 closure parameter that sets up the oscillation type is different, of course.

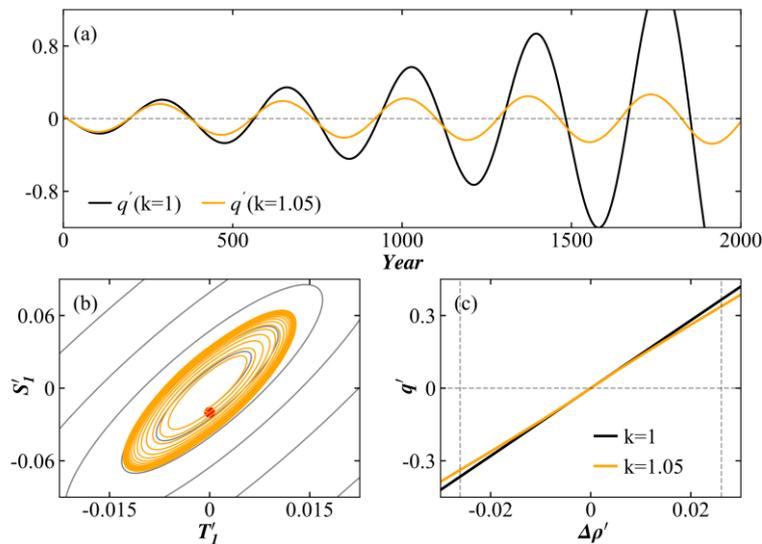
62 Second, we would like to emphasize that to realize a self-sustained oscillation in the linear
63 system, a certain degree of nonlinearity is needed. In the salinity-only system of LY22, enhanced
64 vertical salinity mixing (i.e., nonlinear mixing) was considered in the subpolar ocean. Therefore,

65 under a reasonable linear closure parameter λ , the enhanced vertical salinity mixing will turn the
 66 growing oscillation of the 4S model into a self-sustained oscillation. We have stressed that this
 67 mixing in the subpolar ocean cannot be too strong; otherwise, the growing oscillation will become a
 68 damped oscillation and the 4-box model will practically become the 3-box model. There is no self-
 69 sustained oscillation in the *linear* 3-box and 4-box models. The self-sustained oscillation can only
 70 occur when a certain degree of nonlinearity is considered. It cannot be understood as “*a mix of the*
 71 *unstable oscillation of 4S and the stable one of 3S.*”

72 *I think this is the major problem of the paper. All these things need to be discussed and explain*
 73 *in length. I would like an alternative explanation for the existence of self-sustained oscillation*
 74 *despite the instability of both 3TS and 4TS (for $\lambda=14 \text{ Sv Kg}^{-1} \text{ m}^3$).*

75 We are sorry for not having explained it well in previous manuscript. In the revised manuscript
 76 we state: “**The essence for a self-sustained oscillation is a linearly growing oscillation**
 77 **restrained by a nonlinear process, which can take the form of a nonlinear subpolar vertical**
 78 **mixing, or of a nonlinear relation between AMOC anomaly and meridional difference of**
 79 **density anomaly, or take other nonlinear forms.**”

80 Therefore, even in a 3-box system, if assuming a small degree of internal nonlinearity between
 81 AMOC anomaly and meridional difference of density anomaly (Fig. R1c, orange curve), a self-
 82 sustained oscillation can occur (Figs. R1a, b, orange curves). This is an important development of
 83 LY22. If the 3-box system contains no nonlinear factors, self-sustained oscillation cannot occur
 84 (Figs. R1a, b, black curves). This further suggests that the nonlinearity is the key to self-sustained
 85 oscillation.



86
 87 FIG. R1. Oscillations under $\lambda = 14 \text{ Sv} \cdot \text{kg}^{-1} \text{ m}^3$ for the 3TS model. (a) Time series for q' (units: Sv) under
 88 $k = 1$ (black curve) and $k = 1.05$ (orange curve). (b) T_1' - S_1' phase space diagrams for years 1-10000. The red

89 dot represents the initial location of T'_1 and S'_1 . Black curve is for $k = 1$, and orange curve is for $k = 1.05$.
90 (c) Variation of q' with $\Delta\rho'$ (units: kg/m^3) under $k = 1$ (black curve) and $k = 1.05$ (orange curve). The
91 intersections between the vertical dashed gray lines and the abscissa axis mark the upper and lower limits for
92 $\Delta\rho'$ during the integration. The values of the other parameters are the same as those listed in Table 1 in the
93 revised manuscript.

94

95 *However, in the present study, according to Figure 2c, 3TS and 4TS both exhibit stable*
96 *oscillation under the standard value ($\lambda=12$ Sv Kg-1 m3) and both exhibit unstable oscillation*
97 *under the other tested value ($\lambda=14$ Sv Kg-1 m3). In either case the explanation of LY22*
98 *failed. The addition of mixing in the subpolar region leading to a self-sustained oscillation cannot*
99 *be interpreted as a mix of a stable oscillation in 3TS and an unstable one in 4TS. This makes the*
100 *entire purpose of the theoretical model useless.*

101 Since the temperature equations affect the behaviors of the system, it is natural that the critical
102 values of λ (i.e., λ_c) in the 4S (3S) and 4TS (3TS) models differ. This does not necessarily suggest
103 the failure of the explanation of LY22, since a self-sustained oscillation is not a result of “a mix of a
104 stable oscillation in 3TS and an unstable one in 4TS.” A self-sustained oscillation emerges from a
105 growing oscillation that is restrained by a nonlinear process.

106 The theoretical models used in LY22 and this manuscript are extremely useful. The linear
107 stability analyses on the linear model give us the conditions of stability and the oscillation of the
108 system. Particularly, the 3S model of LY22 can be solved completely theoretically, so that we can
109 see exactly which model parameters and how these parameters determine the stability and the
110 oscillation of the system (see section 4c of LY22). This allows a fundamental understanding of the
111 model behaviors.

112 *Finally, the regime with $\lambda \sim 13$ Sv Kg-1 m3, appear quite interesting since it provides a*
113 *stable oscillation for 4TS and an unstable oscillation for 3TS (Fig.2c). This is the opposite to LY22*
114 *for $\lambda \sim 12$ Sv Kg-1 m3. How would you explain the increase stability of the 4-box model*
115 *over the 3-box model?*

116 We rephrase your above concern as “why the 4TS model is more stable than the 3TS model,
117 while in LY22 the 3S model is more stable than the 4S model, regardless of the value of λ (Fig.
118 2c)?” We have a detailed explanation in lines 388-420 of the revised manuscript, though here we
119 would like to provide a more straightforward explanation in the following context.

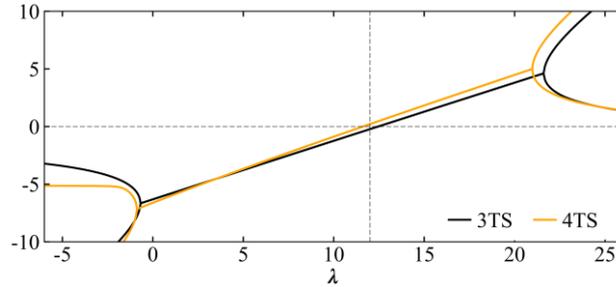
120 Suppose there is an initial perturbation of the AMOC (q'), when $q' > 0$ ($q' < 0$), the subpolar
121 salinity and temperature perturbations $S'_2 > 0, T'_2 > 0$ ($S'_2 < 0, T'_2 < 0$). In the salinity-only 4S
122 system of LY22, $S'_2 > 0$ ($S'_2 < 0$) leads to a stronger (weaker) downward motion in the subpolar

123 ocean, so that it in turn reinforces the initial q' . This process can be simply sketched as $q' > 0 \rightarrow$
 124 $S'_2 > 0 \rightarrow q' > 0$, or $q' < 0 \rightarrow S'_2 < 0 \rightarrow q' < 0$. In other words, the initial q' can be reinforced by
 125 the subpolar vertical salinity perturbation; or the subpolar salinity stratification has potentially a
 126 destabilizing effect on the oscillation. This destabilizing effect is absent in the 3S model of LY22.
 127 Therefore, the 3S model is more stable than the 4S model.

128 In a *temperature-only* box model, we have the following process: $q' > 0 \rightarrow T'_2 > 0 \rightarrow q' < 0$,
 129 or $q' < 0 \rightarrow T'_2 < 0 \rightarrow q' > 0$. That is, the initial q' can be damped by the subpolar vertical
 130 temperature perturbation. In other words, the subpolar temperature stratification has potentially a
 131 stabilizing effect on the oscillation. This stabilizing effect is absent in the 3T model. Therefore, the
 132 4T model is more stable than the 3T model.

133 In the temperature-salinity box model (4TS), the resultant effect of the subpolar temperature
 134 and salinity stratifications on the oscillation behavior is determined the relative effects of the two.
 135 We show in the manuscript that under a reasonable surface temperature restoring timescale, the
 136 stabilizing effect of temperature stratification will overcome the destabilizing effect of salinity
 137 stratification; therefore, the 4TS model is more stable than the 3TS model.

138 This does not necessarily suggest inconsistency between the TS models and the LY22 models.
 139 We can also show that the 3TS model is more stable than the 4TS model, given some unrealistic
 140 surface temperature restoring timescale (Fig. R2).



141
 142 Fig. R2. Dependences of real parts of eigenvalue ω on λ in the 4TS (orange curve) and 3TS models (black
 143 curve) under the temperature restoring coefficient $\gamma = (1 \text{ month})^{-1}$. The units of the ordinate are 10^{-10} s^{-1} .
 144 The values of the other parameters are the same as those listed in Table 1. The vertical dashed line denotes
 145 the situation under the standard value $\lambda = 12 \text{ Sv} \cdot \text{kg}^{-1} \text{m}^3$.

146
 147 **2. Inconsistency with hypothesis of LY22 when Temperature anomaly are allowed**

148 *I think this is a major problem of the paper. What you demonstrated here is that the self-*
 149 *oscillation does not work at all like LY22.*

150 *The addition of T creates a regime where both 4TS and 3TS are both stable (under $\lambda =$*
151 *12 Sv kg-1 m3) or unstable (under $\lambda = 14$ Sv kg-1 m3). This is a huge discrepancy with 3S*
152 *and 4S model - and the explanation of the self-sustained oscillation in LY22.*

153 **Responses:** Thank you very much for these comments. The self-sustained oscillation mechanism in
154 LY22 and that in this work are consistent, and can be concluded as follows: **“a linearly growing**
155 **oscillation restrained by a nonlinear process,”** instead of “a mix of a stable oscillation in 3TS and
156 an unstable one in 4TS.” Here, we would like to emphasize again that the results from the 4TS
157 (3TS) model is mostly consistent with the results from the 4S (3S) model. The only difference is
158 that the 3S model is more stable than the 4S model, while the relative stability of the 3TS and 4TS
159 models is not straightforward, depending on the strength of surface temperature restoring. The
160 addition of temperature will naturally affect the stability of the models, thus the relative stability
161 between 3-box and 4-box models. If there is no difference in system stability after including
162 temperature, our current work “role of temperature” would become useless. Moreover, the more
163 stable 4TS model than the 3TS model suggests exactly that the temperature effects make the 4-box
164 models more stable than the 3-box models.

165 The difference in relative *linear* stability between 3-box and 4-box models is not important,
166 because as far as the self-sustained oscillation is concerned, the fundamental mechanism is the same
167 for all the models; that is, **the self-sustained oscillation emerges from a linearly growing**
168 **oscillation that is restrained by a nonlinear process.”** The nonlinear process can be enhanced
169 vertical mixing in the subpolar ocean, or a nonlinear relationship between AMOC anomaly and the
170 meridional difference of density anomaly, or be other nonlinear forms. Therefore, the self-sustained
171 oscillation mechanism in LY22 and that in the TS models are consistent, and the change in relative
172 stability between 3-box and 4-box models after including temperature is reasonable.

173 *One perfect example is the use of the dashed orange curves in Figs. 8, 9, and 10. There*
174 *definition was clear in LY22 ("end" of the self-oscillatory regime, through the instability of 3S).*
175 *Here it is not explained.... (Note that the orange computation of the solid line is also quite unclear -*
176 *Is it the instability of the 4TS, 3TS, or both?). It looks like a lucky guess. This is not acceptable. This*
177 *is highly problematic. How did you compute these two lines and how do they relate to 4TS and 3TS*
178 *stability? If it does not relate to it, this suggests that the computation of the stability of 4TS and 3TS*
179 *is not useful. And that the system does not behave equivalently to LY22 when T is introduced.*

180 *This raises questions regarding the robustness of LY22. My personal conclusion after reading*
181 *the manuscript is that LY22 results cannot be generalized to the presence of temperature. A detailed*

182 *discussion on that is needed. Qualitative inconsistencies with LY22 should be clarified (i.e.,*
183 *difference between S and TS model).*

184 Sorry for not having clearly described the meaning of the solid and dashed orange curves in
185 Figs. 8-10. The solid and dashed orange curves in both LY22 and this manuscript stand for the
186 lower and upper stability limits within which the self-sustained oscillation can occur, when
187 including an enhanced subpolar vertical mixing. They are all calculated by numerical integration;
188 no self-sustained oscillation occurs when the system is less stable than the dashed orange curve in
189 the 4TS (4S) models, no matter how strong the subpolar vertical mixing is. The stability threshold
190 of the 3S model is also higher than that of the 4S model, and can be well represented by the dashed
191 orange curve.

192 The stability threshold of the 3TS model under the parameters of Figs. 8-10 is lower than that
193 of the 4TS model, as revealed in Fig. 2; thus, it will never become the upper stability limit for the
194 self-sustained oscillation. Therefore, the dashed orange curve in Figs. 8-10 of the original
195 manuscript stands for the upper stability limit for the self-sustained oscillation when subpolar
196 vertical mixing is included, instead of the stability threshold of the 3TS model. Even so, the
197 calculation of 3TS stability is meaningful, since the change of relative stability between the 3-box
198 and 4-box models after including thermal effects can reflect the effects of thermal processes. The
199 more stable 4TS (3TS) model than the corresponding 4S (3S) model also suggests that the thermal
200 effects have stabilizing effects on the system. In this way, the topic “role of temperature” is
201 justified. The system behavior will naturally differs from that of LY22 since the thermal effects
202 must play a role in the system stability, but the underlying self-sustained oscillation mechanisms in
203 LY22 and that in our current work are still consistent. We made it clear in this revised manuscript
204 that the solid orange curve stands for both the stability threshold and the lower stability limit for the
205 self-sustained oscillation in the 4TS model. For simplicity, we removed the dashed orange curve
206 and related discussion.

207 Please also refer to our reply to your 1st question. This paper is a further development of LY22.
208 The self-sustained oscillation mechanism raised in LY22 can be generalized to cases with both
209 temperature and salinity. Figure 2 in the manuscript shows clearly the curves for 3S, 4S, 3TS, and
210 4TS are almost identical, except that the values of λ_C , λ_1 and λ_2 differ in different models. Note that
211 the curves for the 3S model are obtained purely theoretically (LY22), while the curves for the other
212 models are obtained numerically. The consistency between these curves is not of coincidence, but
213 of certainty, because of the consistency of the physical fundamentals in these models. Different
214 models and different physical processes lead to slightly different sensitivity of the AMOC anomaly
215 to the meridional difference of density anomaly.

216

217 3. *Lack of information*

218 *The manuscript suffers from a lack of relevant and clear information. This affects both the*
219 *understanding of the study and its reproducibility. This needs to be fix.*

220 *It should be made clear that you have two models. In my understanding, since you impose the*
221 *mean state, you have an only-salinity *active* variable model and a both TS *active* variable*
222 *model. Where active mean that anomaly can exist (i.e., $\dot{T/S}$). Clarifying this point in the text*
223 *would help the reader.*

224 *If you want to set "any" background state. You should discuss the equilibrium (2a-b) with care,*
225 *since they are not used in full. To my understanding you need to set the same background (and not*
226 *the independent equilibrium of 3TS 4TS, 3S, and 4S) to have a fair comparison between the two*
227 *versions of the model TS and S. Otherwise it is impossible to use the exact same background state*
228 *for each version. Do I understand correctly? If so, could you add this rationale in the text?*

229 *In equation (7b) and (7e) T2, S2, and V2 does not have the same meaning. For instance V2 is a*
230 *larger volume in 7 than in 5. This should be reflected by the use of different symbols.*

231 *Without the explanation or inclusion of the equations for the S model (derived in LY22). It is*
232 *hard to follow. Could you confirm that the S model is actually the (5 e-h) and (7d-f)?*

233 **Responses:** Thank you very much for these suggestions. We clarified in the revised manuscript that
234 we use models with only salinity equations (4S and 3S models) and models with both temperature
235 and salinity equations (4TS and 3TS models).

236 We demonstrated more clearly in this revision that Eqs. (5e-h) and Eqs. (7d-f) are the models
237 of LY22, while in the 4TS and 3TS models both temperature and salinity are included.

238 The equilibrium states in Eqs. (2a-b) have been fully used in Eqs. (5a-h), but might be
239 confusing due to the linearization from Eqs. (1a-h) to Eqs. (5a-h). Take Eqs. (1e) and (1g) as
240 examples, linearizing them gives $q'(\overline{S}_4 - \overline{S}_1) + \overline{q}(S'_4 - S'_1) + F_w$ and $q'(\overline{S}_2 - \overline{S}_3) + \overline{q}(S'_2 - S'_3)$ on
241 the right-hand side, which can be finally reduced to Eqs. (5e) and (5g). Therefore, some equilibrium
242 values like \overline{T}_3 and \overline{S}_3 are cancelled, and all the equilibrium values have been used exactly. The
243 equilibrium values for temperature are calculated after imposing the same pair of restoring
244 temperatures T_1^* and T_2^* for the 4TS and 3TS models. Based on Eq. (2a), the equilibrium values for
245 temperature must be different in the 4TS and 3TS models since the 3TS model has a larger V_2 .
246 Other standard background parameters are the same in the 3TS, 4TS, 3S, and 4S models, while the
247 corresponding rationales have been added in the revised manuscript. In conclusion, we have set the

248 same background states for the 3TS, 4TS, 3S, and 4S models, while the calculated equilibrium
249 temperatures must differ.

250 To be consistent with LY22, we still use V_2 in the 3-box models; we stress that the V_2 in the 3-
251 box models equals to the sum of V_2 and V_3 in the 4-box models.

252

253 4. *Physical description*

254 *I have several issues with the physics described in the text.*

255 **Responses:** Thank you very much for these comments.

256 *I disagree that the feedbacks could be spotted/illustrated by the lag of the timeseries. The lag is*
257 *mainly related to the oscillation which is driven (as well as its timescale) by the mean advection.*

258 We agree that the lag is mainly due to the mean advection. In the manuscript, we state that the
259 positive (negative) correlation coefficient at **lag 0** is a further illustration of the positive (negative)
260 feedback.

261 We would like to emphasize that the feedbacks we discussed in the correlation figures are
262 based on the deterministic equations, not just deduced from the correlations. If there is no equation
263 between different processes, the correlation between different processes suggests exactly some kind
264 of relationship and one can never say there is causality between them.

265 In the correlation figures of the manuscript, all processes are connected by equations.
266 Therefore, we can say that the positive (negative) correlation coefficient at lag 0 is a further
267 illustration of the positive (negative) feedback, because the bases are in the equations, thus the
268 underlying physics, instead of correlation coefficients alone.

269 *A description of what the nonlinear AMOC-density relation physically represents would be*
270 *useful. Also, a description of how it is parameterized in your equation (10) would be useful.*

271 This nonlinear relation stands that if the meridional difference of density anomaly is larger
272 than certain threshold value (like ρ_{cri} in section 4b), the growth of AMOC anomaly could show
273 certain degree of nonlinearity. This parameterization follows the one used in Rivin and Tziperman
274 (1997; RT97). In our work, only very weak nonlinearity ($k=1.05$) is considered, much weaker than
275 that used in RT97 ($k=3$). This also suggests that this kind of nonlinearity is quite efficient at turning
276 a linearly growing oscillation into a self-sustained one.

277 More detailed description is added in the revised manuscript.

278 *I had a hard time understanding what you called the "positive restoring-advection feedback".*

279 *For instance, in l.243-245, the fact that T^2 is decreased at the end suggests that you are*
280 *illustrating a negative feedback. I wonder if what you are referring to is not the action of*
281 *temperature restoring on a density anomaly dominated by salinity and partially compensated by*
282 *temperature. Here the action of thermal restoring will reduce the temperature anomaly hence*
283 *intensifying the density anomaly. This is some kind of positive feedback for the density. However it*
284 *cannot be described using only temperature. Also this is more an oddity than a positive feedback.*
285 *Indeed this mechanism still only leads to a transient increase of the density, but asymptotically it is*
286 *still removing perturbation (because it is driven from the negative feedback induced by the*
287 *temperature surface restoring). This should be clarified*

288 Thank you very much for pointing out that the “positive restoring advection feedback” had not
289 be described well in the original manuscript.

290 In line 237 of the original manuscript, we stated: “There are mainly two feedbacks between the
291 thermal processes and the AMOC.” The positive restoring advection feedback describes a relation
292 between temperature restoring and AMOC (advection), not between temperature restoring and
293 temperature itself.

294 As you point out above, “*the action of thermal restoring will reduce the temperature anomaly*
295 *hence intensifying the density anomaly. This is some kind of positive feedback for the density.*” This
296 is correct since the density anomaly determines AMOC anomaly, and they are positively correlated.
297 It is also true that the density anomaly is not only determined by temperature. Here, mainly
298 temperature-related density anomaly and thus temperature-related AMOC anomaly is discussed.

299 In the revised manuscript, this feedback is discussed in more detail. Take box 2 as an example,
300 by limiting the increase of subpolar temperature anomaly and thus the negative temperature-
301 advection feedback, the restoring effect manifests as the positive feedback between the restoring
302 term and AMOC anomaly. In other words, such restoring advection feedback is to increase AMOC
303 anomaly, i.e., to amplify the initial AMOC perturbation.

304 Note that the temperature feedbacks in the 4TS model should be viewed as a combination of a
305 negative temperature advection feedback and a positive restoring advection feedback. The latter is
306 driven by the former, which in turn hampers the former. Their total effect is a negative feedback.
307 Nevertheless, the restoring advection feedback should still be termed as a positive one.

308 *Mean advection described as a feedback can be misleading. This is not exactly the same kind of*
309 *feedback as the positive salinity feedback, for instance. Mean advection will conserve the anomaly*
310 *(even in a linear framework), it just moves things around. It leads to an oscillation of period*

311 $\sim 1/\bar{q}$. This can be seen by the action of the mean advection in the equations for the anomaly:
312 it acts as a skew-symmetric component of the Jacobian operator.

313 We totally agree with your comments. Figure 3b shows that the mean advection term has the
314 smallest contribution to the temperature anomaly, although this term has good positive correlation
315 with temperature-related AMOC anomaly (Fig. 3e). Therefore, practically the positive mean
316 advection feedback can be neglected.

317 In the revised manuscript, we tone down the mean advection feedback to avoid misleading our
318 readers.

319

320 5. Choice of parameters

321 *The choice of the parameter is not well discussed.*

322 **Responses:** Sorry about this.

323 *For instance, the choice between $\lambda = 12 \text{ Sv } \cdot \text{kg}^{-1} \text{m}^3$ (as in LY22) or $\lambda = 14 \text{ Sv } \cdot \text{kg}^{-1}$
324 m^3 here is not discussed. Also the value of $13 \text{ Sv } \cdot \text{kg}^{-1} \text{m}^3$ seems more interesting since it provide a
325 regime where 3TS and 4TS are unstable and stable, respectively. All these choices need to be
326 explained.*

327 In the revised manuscript, more explanation on the choice of λ is provided.

328 The choice of $\lambda = 12 \text{ Sv} \cdot \text{kg}^{-1} \text{m}^3$ as the standard parameter is to make it the same as in
329 LY22. Under $\lambda = 12 \text{ Sv} \cdot \text{kg}^{-1} \text{m}^3$, the 4S model is unstable while the 4TS model is stable,
330 **reflecting that the overall thermal effect is to stabilize the system.**

331 The choice of $\lambda = 14 \text{ Sv} \cdot \text{kg}^{-1} \text{m}^3$ is to make the 4TS model unstable, since the self-sustained
332 oscillation has to be based on a linearly unstable regime.

333 Because the subpolar temperature stratification effect is a stabilizing effect and can overcome
334 the destabilizing effect of subpolar salinity stratification under realistic parameter ranges, the 4TS
335 model can be more stable than the 3TS model. Therefore, the choice of $\lambda = 13 \text{ Sv} \cdot \text{kg}^{-1} \text{m}^3$ is to
336 make the 4TS model stable and the 3TS model unstable. However, we cannot realize self-sustained
337 oscillation in the 4TS model under $\lambda = 13 \text{ Sv} \cdot \text{kg}^{-1} \text{m}^3$ since the system stays in a linearly stable
338 regime.

339 *Also the value of $\bar{q} = 10 \text{ Sv}$ is quite low. (GT95 is more consistent to observations.) I
340 would use a value way closer to 20 Sv. Maybe 18 Sv? Do you have any argument to do otherwise?*

341 The choice of \bar{q} is critical to the oscillation timescale. In the observation, the maximum mean
342 AMOC is about 20 Sv. This mass transport includes water in the upper 1000 m and the water from
343 the Southern Hemisphere. In a one-hemispheric box model with the upper ocean depth set to 500 m,
344 the mean mass transport should be remarkably smaller than the realistic value. Otherwise, the
345 turnover timescale for a one-hemisphere box model would be unrealistically short. In Griffies and
346 Tziperman 1995 (GT95), \bar{q} is much larger than that in ours, and their subpolar boxes are much
347 smaller than ours, so that the dominant timescale in the GT95 box model is the decadal timescale,
348 instead of the centennial timescale.

349 Our single-hemispheric model incorporates only the AMOC recirculating in the Northern
350 Hemisphere, so that a smaller mean AMOC is reasonable. This is also consistent with the choice of
351 mean AMOC in Nakamura et al. (1994).

352

353 **Reference:**

354 Nakamura, M., P. H. Stone, and J. Marotzke, 1994: Destabilization of the thermohaline circulation by
355 atmospheric eddy transports. *J. Climate*, 7, 1870-1882.

356

357 *1 year for $1/\tau$ is quite long. It is set to be 60 days for 50 m in NEMO handbook (i.e., $-40 W$*
358 *m⁻² K⁻¹), for example.*

359 We admit that it has been prevailing to set the temperature restoring timescale for a surface
360 layer with a few tens of meters to be 1-2 months in models with higher complexity (Marotzke and
361 Willebrand 1991; Weaver and Sarachik 1991; Mysak et al. 1993; Pierce 1996). However, there is
362 no such thin surface layer in our theoretical model, due to its simplicity. As a substitute, we permit
363 the temperature restoring to happen over the entire depth range of the upper boxes (0-500 m). Such
364 thick surface layer clearly necessitates a much longer restoring timescale. The rather deep restoring
365 depth and long restoring timescale are common in theoretical studies (GT95; Roebber 1995; RT97;
366 Scott et al. 1999; Lucarini and Stone 2005).

367

368 **References:**

369 Griffies, S. M., and E. Tziperman, 1995: A linear thermohaline oscillator driven by stochastic atmospheric
370 forcing. *J. Climate*, 8, 2440-2453.

371 Lucarini, V., & Stone, P. H. (2005). Thermohaline circulation stability: A box model study. Part I:
372 Uncoupled model. *J. Climate*, 18(4), 501-513.

- 373 Marotzke, J., and J. Willebrand, 1991: Multiple equilibria of the global thermohaline circulation. *J. Phys.*
374 *Oceanogr.*, 21, 1372-1385.
- 375 Mysak, L. A., T. F. Stocker, and F. Huang, 1993: Century-scale variability in a randomly forced, two-
376 dimensional thermohaline ocean circulation model. *Climate Dyn.*, 8, 103-116.
- 377 Pierce, D. W., 1996: Reducing phase and amplitude errors in restoring boundary conditions. *J. Phys.*
378 *Oceanogr.*, 26, 1552-1560.
- 379 Rivin, I., and E. Tziperman, 1997: Linear versus self-sustained interdecadal thermohaline variability in a
380 coupled box model. *J. Phys. Oceanogr.*, 27, 1216-1232.
- 381 Roebber, P. J. (1995). Climate variability in a low-order coupled atmosphere-ocean model. *Tellus Ser. A-*
382 *Dyn. Meteorol. Oceanol.*, 47(4), 473-494.
- 383 Scott, J. R., Marotzke, J., & Stone, P. H. (1999). Interhemispheric thermohaline circulation in a coupled box
384 model. *J. Phys. Oceanogr.*, 29(3), 351-365.
- 385 Weaver, A. J., and E. S. Sarachik, 1991: Evidence for decadal variability in an ocean general-circulation
386 model - An advective mechanism. *Atmos. Ocean*, **29**, 197-231.

387

388 6. *Introduction*

389 *The context within the literature is unclear and extremely short. The multi-centennial variability*
390 *of the AMOC in model, theory, and observations is a wide research topic. I expect a more in-depth*
391 *discussion of our current knowledge and of gaps in our current understanding. Introduction should*
392 *be re-written with clear scientific questions and clear discussion of the large literature on the topic.*

393 **Responses:** Thank you very much for this suggestion. The introduction has been rewritten. We
394 reviewed the mechanisms for multicentennial AMOC oscillation in state-of-the-art coupled models,
395 and realized that their discrepancy calls for theoretical studies. On this account, we further reviewed
396 the related theoretical studies, and found that they each reflects a certain degree of inadequacy in
397 accounting for the sustainable multicentennial AMOC oscillation. Finally, we stated the prospective
398 effects that thermal processes have on AMOC oscillation, and came up with our improved model of
399 LY22, namely, temperature variation is permitted now.

400

401 **Specific Comments:**

402 1. *l.31: Change "the thermal process exerts mainly" to "the thermal processes exert mainly".*

403 Revised.

404 2. *l.33: Change "stratification, which are contributed by" to "subpolar stratification. These*
405 *thermal processes are composed of"*

406 **Responses:** Thank you very much for this comment. The abstract is revised; and this sentence is
407 rewritten. We have made it clear: the three thermal effects (shortening the oscillation period,
408 stabilizing the overall system, and stabilizing the subpolar stratification) are separately caused by
409 three thermal processes (the fast surface temperature restoring process, the negative temperature
410 advection feedback, and the stabilizing subpolar temperature stratification).

411
412 3. *l.34: Change "feedback and subpolar" to "feedback, and the subpolar"*

413 **Responses:** Thank you very much for this suggestion. This sentence is rewritten.

414
415 4. *l.35: Remove ", respectively"*

416 **Responses:** Thank you very much for this suggestion. This sentence is rewritten.

417
418 5. *l.54-55: It would be more accurate to state that only the anomalous salinity evolution was kept.*
419 *Indeed in these studies the temperature was acknowledged within the background state (making*
420 *them more "realistic"). Unlike the salt-oscillator of Huang and Dewar (Journal of Physical*
421 *Oceanography, 1996), for instance.*

422 **Responses:** Thank you very much for this suggestion. The introduction has been rewritten; and this
423 sentence is deleted. In the revised manuscript, we clarified that in LY22 model only salinity
424 variation (or evolution) is possible while in our current model both salinity and temperature
425 variations (evolutions) are permitted.

426
427 6. *l.65-67: Describe like that it is a negative feedback. In general restoring is a negative feedback*
428 *(acting along the diagonal of the Jacobian Matrix). Please clarify.*

429 **Responses:** Thank you very much for this comment. We have rewritten the introduction; and here
430 we give a clarification on this problem. The feedback between restoring and temperature is indeed
431 negative, which could be termed as the negative restoring temperature feedback, while the restoring
432 advection feedback (which we actually focus on in the paper) depicts the relation between restoring
433 and AMOC perturbation.

434

435 7. l.90: add a coma after "properties". Revised.

436 8. l.98: Change "Model formulae" to "Model formulation". Revised.

437 9. l.107-108: A bit odd to use τ (which is often use for time "t") as an inverse of a time scale. I
438 suggest to use γ , for instance, or to write as $1/\tau$.

439 **Responses:** Thank you very much for this suggestion. In the revised manuscript, we replaced τ with
440 γ everywhere.

441

442 10. l.122: remove "diagrams". Revised.

443 11. l.124: Change "lower tropical oceans" to "deeper tropical ocean boxes". Revised.

444 12. l.124-125: Change "lower subpolar oceans" to "deeper subpolar ocean boxes". Revised.

445 13. l.125: Change "lower ocean depths" to "deeper ocean box depths". Revised.

446 14. l.129: remove "easily". Revised.

447 15. l.157: add a coma after " α ". Revised.

448 16. l.158: add a coma after "expansion". Revised.

449 17. l.217: Would "increase frequency" be better than "acceleration of the oscillation"?

450 **Responses:** Thank you very much for this suggestion. We rephrased the sentence accordingly.

451

452 18. l.219: Change "for" to "of". Revised.

453 19. l.275: Change "very front part" to "initial part"

454 **Responses:** Thank you very much for this suggestion. This sentence is revised; and the paragraph
455 containing this sentence is moved to the end of section 2a.

456

457 20. l.279: It would be more interesting to (also) discuss the ratio of τ over \bar{q} , rather than
458 τ alone. I.e., It would be more physical to compare the relative action of two processes.

459 **Responses:** Thank you very much for this comment. The primary object of this study is the role of
460 temperature in multicentennial AMOC oscillation. Thus, we first discussed the effects of
461 temperature feedbacks (of course including the restoring advection feedback whose strength is

462 determined by γ , the substitute for τ in this version of manuscript) in section 3. As for \bar{q} (related to
463 flow property), we put it in the part of sensitivity studies (section 5). Additionally, we treated the
464 model parameters as being independent of each other in our model; thus, we tested eigenmode's
465 sensitivity to each parameter, instead of to multiple parameters like γ/\bar{q} .

466

467 *21. l.296-297: This is because the restoring for tau->0 acts essentially as a flux. So you are back in*
468 *a system equivalent to the salinity-only one. (I.e., flux boundary condition for both T and S and*
469 *not mixed boundary condition.) In this context everything could be written with a single*
470 *variable (i.e., density).*

471 **Responses:** Thank you very much for this comment. Under $\gamma = 0$, $\bar{T}_1 = \bar{T}_2 = \bar{T}_3 = \bar{T}_4$; thus, the
472 temperature advection becomes null. Therefore, the density variation will be exclusively controlled
473 by salinity variation at this point.

474

475 *22. l.299: It is even stronger. You have: T1=T*1 and T2=T*2.*

476 **Responses:** Thank you very much for this comment. We clarified it in the revised manuscript.

477

478 *23. l.337-341: I am not sure to follow... Are you suggesting that the action of restoring (what you*
479 *call "they are removed") is acting more efficiently in 4TS than 3TS? I wonder if this does not*
480 *come from the fact you use the same τ for both 4TS and 3TS. This is quite unphysical that a*
481 *restoring will act as quickly on a 500 m layer and on a 4000 m layer. Maybe τ should be a*
482 *function of the thickness of the layer - to be closer to a constant flux (as often hypothesized in*
483 *more advanced numerical model: $-40 \text{ W m}^{-2} \text{ K}^{-1}$, as suggested in the NEMO handbook, for*
484 *instance). Overall there is a problem in the lack of equations/explanations of the 3TS model.*
485 *One have to guess your treatment of box2 in 3TS...*

486 **Responses:** Thank you very much for this comment. Our explanation of the 3TS model was not
487 clear, which led to such misunderstanding. The “removed” actually denotes the advection of
488 anomalies out of the subpolar region. Since there are no temperature and salinity differences among
489 boxes 2, 3, and 4, it is the mean advection of anomalies that transports the anomalies out of the
490 subpolar region in the 3TS model. Since boxes 2 and 3 are well mixed, the transportation time of
491 anomalies between boxes 2 and 3 is saved.

492

493 24. l.347-361: *I find it hard to follow, with discussion of stratification of the subpolar (which does*
494 *not exist in 3S or 3TS) and no figure of its evolution has a function of τ . The lag correlation*
495 *is not a demonstration of a stabilizing or destabilizing effect. I do not understand this argument.*
496 *(Lag-)Statistical link (which are not causality) cannot demonstrate a dynamical feedback.*
497 *Overall I think the problem is simpler: restoring -> stabilizing effect. Their different impact on*
498 *3TS and 4TS is not, for me, demonstrated here by this discussion.*

499 **Responses:** Thank you very much for this comment. We have rewritten these sentences.

500 The correlation in these figures can be a demonstration of a stabilizing or destabilizing effect,
501 because all processes are connected by deterministic equations. If there is no equation connecting
502 different processes, the correlation between different processes can only suggest some kind of
503 relationship, and one can never say there is causality between them.

504 In the manuscript, we state that the positive (negative) correlation coefficient at **lag 0** is a
505 further illustration of the positive (negative) feedback, because the bases are the equations and the
506 underlying physics. Therefore, the statistical link in this work can demonstrate a dynamic feedback.

507 Compared to the 3TS model, the subpolar stratification in the 4TS model ($T'_2 - T'_3$) increases
508 the time that the subpolar temperature anomaly stays within the subpolar region; thus, it physically
509 has a stabilizing effect. This process is lacking in the 3TS model. After this physical interpretation,
510 we can use the correlation coefficient at lag 0 as a more intuitional representation for the
511 stabilizing/destabilizing effect.

512

513 25. l.386-388: *It would be nice to add something like: "If strong enough, the mixing make the 4TS*
514 *model virtually equivalent to the 3TS one".* Revised.

515 26. l.400: *You need to show longer simulations. It is unclear for me if the cycle has indeed*
516 *saturated. The figure suggests that there is a slowdown in the rate of growth but not a stability*
517 *in the amplitude of the cycle, yet. A large number of period with the same oscillation amplitude*
518 *would be needed to demonstrate that.* Revised.

519 27. l.411: *Change "insensitive" to "almost insensitive".* This sentence is deleted.

520 28. l.412: *Change "dominated" to "controlled".* This sentence is deleted.

521 29. l.414-415: *Change "lead to" by "stabilize the unstable oscillation and lead to".* This sentence
522 is rewritten.

523 30. l.431-432: *Yes, but the parameter also "are hardly deviated". So I don't think it demonstrates*
524 *the "robustness". It rather shows that the stability does affect the period. Please rephrase.*

525 **Responses:** Thank you very much for this comment. The text on the impact of mixing and the
526 nonlinear relation between AMOC anomaly-meridional difference of density anomaly on the period
527 was removed.

528

529 31. l.447-451: *You should be slightly more careful and clear. Your background state is set with*
530 *parameters, rather than computed as an equilibrium (as in the previously cited study). Hence,*
531 *despite some advantages (as setting the parameters to what you wish), you take the risk of using*
532 *un-consistent parameters. This should be highlighted.*

533 **Responses:** Thank you very much for this suggestion. This paragraph is deleted. We highlighted
534 that our background state is set with parameters.

535

536 32. l.453: *I do not find it extremely useful to compare with so much details previous (rather old)*
537 *studies. I feel that idealized model are useful to explain physics, rather than to validate each*
538 *other.*

539 **Responses:** Thank you very much for this suggestion. We removed the annotations related to GT95
540 and RT97 in Fig. 8. However, it would be nice to see that different theoretical works are consistent.

541

542 33. Fig8-9-10: *You show us that the period and e-folding timescale can be a bit affected by the*
543 *choice of 3 or 4 boxes (Fig. 4). It would need to be clear which one you choose for the plots. A*
544 *comment on this issue (since GT95 and RT97 are un-consistent on this matter) would be useful.*

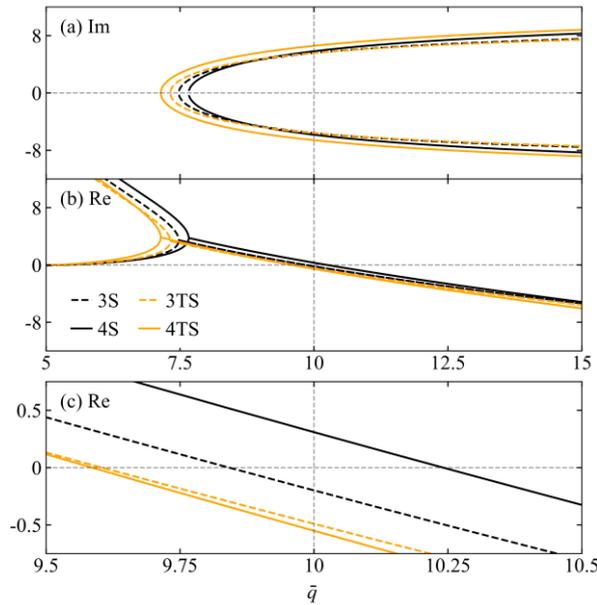
545 **Responses:** Thank you very much for this suggestion. Figures 8-10 are all based on the 4TS model.
546 We clarified this issue in the revised manuscript.

547

548 34. l.499-500: *Technically if a term is used to build the period (\bar{q}) it could not affect the*
549 *growth. This could be easily shown by computing the eigenvalues by hand. Terms that end up in*
550 *the imaginary part cannot contribute anymore to the real part of the eigenvalues.*

551 **Responses:** Thank you very much for this comment.

552 We plotted the dependence of real and imaginary parts of the eigenmode on \bar{q} in Fig. R3 under
 553 the standard parameters in Table 1. We find that \bar{q} influences both the oscillation period and e-
 554 folding time of the system. \bar{q} affects the system stability through the mean advection feedback, and
 555 affects the oscillation period through its influence of the overturning rate. Overall, although we are
 556 unable to solve the analytical solution of the 4TS model theoretically, \bar{q} still influences both the e-
 557 folding time and period of the eigenmode, as revealed by linear stability analysis (Fig. R3).



558
 559 FIG. R3. Dependences of (a) imaginary parts and (b) real parts of eigenvalue ω on \bar{q} in the 4TS (solid orange
 560 curves), 3TS (dashed orange curves), 4S (solid black curve), and 3S (dashed black curves) models under $\lambda =$
 561 $12 Sv \cdot kg^{-1}m^3$. (c) is the zoomed-in version of (b) near line $Re(\omega) = 0$. Results of the 4S and 3S models
 562 are from LY22. The units of the ordinate are $10^{-10} s^{-1}$. The values of the other parameters are the same as
 563 those listed in Table 1. The vertical dashed line denotes the situation under the standard value $\bar{q} = 10 Sv$.

564
 565 *35. Fig.9: It is problematic that the realistic regime $\bar{q} > 14$ corresponds to regime 1. A word*
 566 *on that would be useful.*

567 **Responses:** Thank you very much for this comment. We have removed the annotations of regime 3.
 568 We think that under realistic \bar{q} , the system should be unstable and there is possibility for self-
 569 sustained oscillation. The problem is that in our theoretical model, the depth for boxes 1 and 2
 570 represents the upper ocean instead of the entire AMOC upper branch, and only the AMOC
 571 recirculating in the Northern Hemisphere is considered. Thus, the realistic value for \bar{q} should be a
 572 smaller value like 10 Sv, instead of a much larger value like >14 Sv. The standard and realistic $\bar{q} =$
 573 $10 Sv$ corresponds to an unstable regime; and self-sustained oscillation has possibility to occur;
 574 thus, is reasonable.

575

576 *36. section 5c: This is a perfect example why I feel that you should stress that background mean*
577 *state is treated as independent parameters: here an increase in the restoring temperature*
578 *difference does not affect the mean flow (whereas in general it should).*

579 **Responses:** Thank you very much for this suggestion. In the revised manuscript, we clarified that
580 the background parameters are all independent of each other.

581

582 *37. l.523: Change "gradient" to "difference"* Revised.

583 *38. l.530-533: Unclear, please rephrase.*

584 **Responses:** Thank you very much for this comment. We have rewritten the abstract and also
585 revised section 3 of the manuscript, emphasizing that three thermal effects are introduced when
586 considering the temperature equations: (1) an increase of the oscillation frequency, (2) a
587 stabilization of the overall system, and (3) a stabilization of the subpolar stratification. Additionally,
588 these three thermal effects are, respectively, caused by three thermal processes: (1) the fast surface
589 temperature restoring, (2) the negative temperature advection feedback, and (3) the stabilizing
590 subpolar temperature stratification.

591

592 *39. l.522: "consume" I do not understand this word in this context. Please rephrase.* This sentence
593 is deleted.

594 *40. l.552: Change "system" to "oscillation".* This sentence is deleted.

595 *41. l.558-559: There is even more recent studies in 2 other CMIP6 models which show this kind of*
596 *multi-centennial oscillation. This should be mentioned in the introduction.*

597 <https://doi.org/10.1007/s00382-022-06534-4>

598 <https://doi.org/10.1029/2020MS002366>

599 **Responses:** Thank you very much for this suggestion. We rewrote the introduction, and cited these
600 papers.

601

602 *42. l.562: Change "easy to realize" to "included in some form"* Revised.

603 *43. l.569-573: I am not sure to follow the arguments... It is either quite general (and I don't see the*
604 *point) or the relation to your study is not well explained.* This sentence is deleted.

605 44. l.576: Change "*On one hand*" to "*On the one hand*" We rewrote this sentence.

606 45. l.578-579: Change "*as revealed in this paper*" to "*consistently with our study*" We rewrote this
607 sentence.

608 46. l.579-580: Unclear. Please rephrase.

609 **Responses:** Thank you very much for this comment. We revised this sentence as follows: “this also
610 implies that the period of the multicentennial AMOC oscillation is likely to be lengthened in the
611 future, which has not gained attention yet.”

612

613 47. l.581: *I don't know what you mean there with "portion"*. We rewrote this sentence.

614 48. l.590: Change "*reserved to*" to "*included*" Revised.

615 49. l.594: Remove "*authenticity of*" Revised.

616

617 **Replies to Reviewer #2:**

618

619 Thank you very much for all of your constructive comments. We have carefully revised our
620 manuscript based on the comments from you and other reviewers. The following are our point-to-
621 point replies.

622 *This manuscript is an extension of a paper by the same authors by considering the effects of*
623 *temperature, and its advection and restoring feedbacks on the self-sustained multicentennial*
624 *oscillation of the overturning circulation. The authors used 4-box, 3-box and diffusive 4-box models*
625 *to show the condition for the existence of such oscillations and showed how temperature made the*
626 *system more stable than the previous salinity-only oscillation. I found the results interesting and*
627 *helpful for us to further understand the behavior of the MOC, especially in the context of ongoing*
628 *climate change. However, I do have some confusions and concerns about the manuscript as it*
629 *currently is. The general major comments are below, followed by minor comments on details.*
630 *Comments are made in the general order of appearance in the manuscript instead of importance.*

631 **Responses:** Thank you very much for your invaluable suggestions, which help us improve the
632 manuscript tremendously. Combining the comments from all the reviewers, we have revised the
633 manuscript primarily in these following aspects:

- 634 1) We have completely rewritten the introduction. Coupled modelling studies on multicentennial
635 AMOC oscillation are synthetically reviewed, the inconsistency among their mechanisms and
636 the necessity for theoretical studies are disclosed. Inadequacy of previous theoretical models in
637 accounting for sustainable multicentennial AMOC oscillation is also discussed. Finally, the
638 potential impacts of thermal processes on AMOC oscillation are raised, justifying the inclusion
639 of temperature effects in this study.
- 640 2) In section 2, the choices of parameters are discussed in more detail.
- 641 3) In section 3, we categorize the thermal effects more precisely. We propose that there are mainly
642 three effects when including the temperature equations: (1) increase of the oscillation frequency,
643 (2) stabilization of the overall system, and (3) stabilization of the subpolar stratification. These
644 three effects are attributed to the following three processes, respectively: (1) fast surface
645 temperature restoring, (2) negative temperature advection feedback, and (3) stabilizing subpolar
646 temperature stratification. Now, it is easier to understand that the behaviors of the temperature-
647 salinity system are different from the salinity-only model in LY22.

648 4) In section 4, we more clearly describe the self-sustained oscillation mechanism. In LY22, we
649 proposed that the nonlinear subpolar vertical mixing is crucial for self-sustained oscillation. In
650 the revised manuscript, we further propose that assuming a nonlinear relation between AMOC
651 anomaly and meridional difference of density anomaly can also lead to self-sustained
652 oscillation. We further show that the self-sustained oscillation mechanism not only agrees with
653 that of LY22, but also advances the theory of LY22.

654 5) Most figures are re-plotted.

655

656 **Major comments**

657 *1. This manuscript shows several versions of box models in which temperature, salinity and*
658 *overturning circulation interact. For me, it seems that the 4-box and 3-box models are two*
659 *limits of the diffusive 4-box model introduced the latest, and the self-sustained oscillation, which*
660 *is the main focus of this manuscript, is realized in the diffusive model. It will be clearer if the*
661 *diffusive model is introduced first and the manuscript can discuss the focus (self-sustained*
662 *oscillation) earlier. Also, it seems to me that the 3-box model is a convective version of the 4-*
663 *box model. However, it makes more sense if the authors set a condition for convection in your*
664 *numerical simulation (say, density of box 2 is larger or equal to density in box 3) for a*
665 *transition between the 4-box and 3-box model, instead of treating them completely separately,*
666 *at least in the numerical part.*

667 **Responses:** Thank you very much for this comment.

668 Both our previous paper (LY22) and this manuscript focus on self-sustained oscillation. In
669 LY22, the self-sustained oscillation was introduced before the 3-box models. In this manuscript we
670 investigate the role of temperature in self-sustained multicentennial AMOC oscillation. We first
671 analyze the role of temperature via *linear* stability analysis. After we recognized the overall
672 stabilizing role of temperature, we naturally came up with the question: will the self-sustained
673 oscillation occur in the linear 4TS system with the stabilizing effect of temperature?

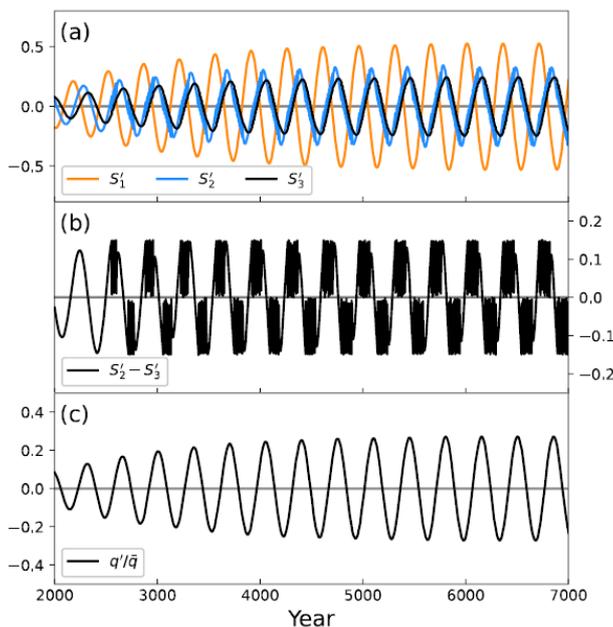
674 The 3-box models are indeed convective versions of the 4-box models. We use the 3TS model
675 to analytically show that the subpolar temperature stratification has a stabilizing effect and can
676 overcome the destabilizing effect of subpolar salinity stratification, in order to reveal the
677 mechanisms of thermal effects and self-sustained oscillation.

678 Thanks for your suggestion of “including convection in numerical simulation to see the
 679 transition between the 4-box and 3-box models.” Note that the convection can only be included in
 680 the 4-box models. The 3-box models are the results after the convection.

681 In an earlier draft of LY22, we did set a condition for convection in the 4S model (similar to
 682 your suggestion) in the numerical simulation; and self-sustained oscillation occurred (Fig. R4). The
 683 convection processes considered included those due to both static instability and convective
 684 instability, which can be expressed as follows,

$$685 \quad |S_2 - S_3| > S_b \Rightarrow \begin{cases} S_2 - S_3 > S_b^+, & \text{Static instability} \\ S_2 - S_3 < S_b^-, & \text{Convective instability} \end{cases}$$

686 where S_b^+ and S_b^- are the thresholds of vertical salinity contrast for convection, whose absolute
 687 values are set to the same for simplicity. Note that S_b^+ and S_b^- can be different, which does not
 688 affect the conclusions of this study. In reality or in ocean-alone and coupled models, the
 689 stratification at the North Atlantic deep-water formation region is such that fresh and cold water lies
 690 on top of saline and warm water (The Lab Sea Group 1998), especially in winter. Under this
 691 background stratification, deep convection can occur.



692
 693 FIG. R4 Self-sustained stable oscillations in the 4S model with both types of deep convection considered. (a)
 694 Time series for S_1' , S_2' , and S_3' (units: psu); (b) time series for $S_2' - S_3'$ (units: psu); (c) time series for q'/\bar{q} .
 695 Before the convection occurs, the evolution of the system in the first 2500 years shows growing oscillation.
 696 In (b), the packed black curves show the happening of deep convection.

697

698 In the final version of LY22 and in this study, we decided enhanced vertical mixing is
699 physically more reasonable than the “one-step” convection in the numerical integration. At least,
700 the changes in salinity and temperature in the subpolar ocean will be much smoother. The
701 convection is just one situation of the enhanced vertical mixing.

702

703 *2. Also, the authors discuss eigenmode sensitivity in section 5, while some sensitivities have*
704 *already been discussed before (e.g., the effect of λ). The structure of the manuscript will be*
705 *clearer if such discussions can be combined.*

706 **Responses:** Thank you very much for this suggestion. The discussion of eigenmode on λ in section
707 2 is to reveal the thermal effects (accelerating the oscillation, stabilizing the overall system, and
708 stabilizing the subpolar stratification). This led to the investigation of thermal processes in section
709 3. To reveal the sensitivity of eigenmode to flow properties, we analyzed the sensitivities of 4TS
710 eigenmode to λ and \bar{q} in section 5b. These two parameters are inseparable since q' is jointly
711 determined by λ and \bar{q} . That is why the discussion of λ occurred twice.

712

713 *3. In the discussion of the 4-box and 3-box models, my impression with the results is that the effect*
714 *of temperature can be described as “marginal” at best (change of period and its contribution in*
715 *q'), even though it indeed makes the system more stable. This might be due to the fact that the*
716 *four- and three-box models are introduced before the diffusive version in which self-sustained*
717 *oscillation is finally emphasized. The manuscript as it is now leaves me the impression that*
718 *temperature is not a crucial factor, even though including temperature is more realistic.*

719 **Responses:** Thank you very much for this comment. Yes, the temperature effects are marginal.
720 However, since in reality there is always temperature, we have to investigate thoroughly the role of
721 temperature. After seeing “temperature is marginal” clearly, we can safely neglect the temperature
722 effects and be more focus on the salinity effects in future studies on this issue.

723

724 **Minor comments**

725 *1. line 33: “caused by” instead of “contributed by”.* Revised.

726 *2. line 38: remove “can”.* We rewrote this sentence.

727 *3. line 39: word “thus”, the causality is not clear in this sentence.* We rewrote this sentence.

728 *4. line 40: “which is less stable” instead of “with less stability”.* We rewrote this sentence.

- 729 5. *line 49: “its thermohaline circulation portion” does not need to include “circulation”.* The
730 introduction has been rewritten.
- 731 6. *line 58: remove “following the Newtonian law”.* The introduction has been rewritten; and this
732 sentence is deleted.
- 733 7. *line 65–67: the description of how restoring-advection feedback is not clear. Do you mean lead*
734 *to more artificial heat loss due to restoring after a positive MOC perturbation? Will this*
735 *actually make the SST lower such that it will hinder the AMOC from recovering or just partially*
736 *hinder warming effect? If there is still warming, AMOC should decrease anyway, maybe to a*
737 *different extent. Some references may be helpful.*

738 **Responses:** Thank you very much for this comment. As you pointed out, more heat will be lost due
739 to restoring after positive MOC perturbation, which was illustrated in Zhang et al. (1993) and
740 Lucarini and Stone (2005). Therefore, the net effect of restoring itself is to hinder the warming
741 effect, although as a whole there is still warming and the AMOC will decrease in strength.

742 The positive restoring advection feedback describes a relation between temperature restoring
743 and AMOC (advection), not between temperature restoring and temperature itself. In the revised
744 manuscript, this feedback is stated clearly with more detail. Take box 2 as an example, by limiting
745 the increase of subpolar temperature anomaly and thus the negative temperature-advection
746 feedback, the restoring effect manifests as a positive feedback between the restoring term and
747 AMOC anomaly. In other words, such restoring advection feedback is to increase AMOC anomaly,
748 i.e., to amplify the initial AMOC perturbation.

749 We would like to emphasize that the temperature feedbacks in the 4TS model should be
750 viewed as a combination of negative temperature advection feedback and positive restoring
751 advection feedback. The latter is driven by the former, which in turn hampers the former. Their
752 combined effect is negative feedback. Nevertheless, the restoring advection feedback should still be
753 termed as positive feedback.

754

755 **References:**

756 Lucarini, V., & Stone, P. H. (2005). Thermohaline circulation stability: A box model study. Part I:
757 Uncoupled model. *J. Climate*, 18(4), 501-513.

758 Zhang, S., R. J. Greatbatch, and C. A. Lin, 1993: A reexamination of the polar halocline catastrophe and
759 implications for coupled ocean atmosphere modeling. *J. Phys. Oceanogr.*, 23, 287-299.

760

761 8. *line 68–69: references needed to verify that the restoring feedback never fully overruns the*
762 *other feedback.*

763 **Responses:** Thank you very much for this comment. We stated that the restoring advection
764 feedback never overruns the temperature advection feedback, instead of the other feedbacks. The
765 temperature feedbacks in the 4TS model should be viewed as a combination of negative
766 temperature advection feedback and positive restoring advection feedback. The latter is driven by
767 the former, which in turn hampers the former. However, their combined effect is negative feedback.

768 Physically, the most extreme form of restoring feedback is that it completely fixes the variation
769 of temperature thus making the total effect of temperature equations null. Therefore, it cannot
770 overrun the temperature advection feedback since the total effect of temperature equations would
771 not be destabilizing. In Nakamura et al. (1994) and Marotzke (1996), the authors described the
772 feedback between meridional atmospheric heat transport and AMOC, which has a similar effect of
773 restoring advection feedback in ocean models employing mixed boundary conditions. It works as
774 follows: an initial positive perturbation of the AMOC lowers meridional temperature difference,
775 weakens meridional atmospheric heat transport, thus lowering high-latitude temperature; the end
776 result is a net strengthening of the AMOC via this feedback alone. Marotzke (1996) further stated
777 that: “since the change in atmospheric transports is a response to anomalous meridional temperature
778 contrasts, which at most can eliminate its cause completely but can never overshoot, this feedback
779 can’t be stronger than the temperature advection feedback.” This supports that the restoring
780 advection feedback cannot overcome the temperature advection feedback.

781

782 **References:**

783 Marotzke, J., 1996: Analysis of thermohaline feedbacks. Decadal Climate Variability: Dynamics and
784 Predictability, D. L. T. Anderson and J. Willebrand, Eds., Springer, 333-378.

785 Nakamura, M., P. H. Stone, and J. Marotzke, 1994: Destabilization of the thermohaline circulation by
786 atmospheric eddy transports. J. Climate, 7, 1870-1882.

787

788 9. *line 75: the comma before the word “whose” is not needed.* The introduction has been
789 rewritten.

790 10. *line 76–79: How do you relate period to the degree of stability of the system? How stable the*
791 *system is determined by the real part of the eigenvalue while period is the imaginary part. And*
792 *some people may intuitively think an oscillation of longer timescale is more stable. Stability and*

793 *period are often discussed together throughout this manuscript, but their relation (if any) seems*
794 *confusing to me.*

795 **Responses:** Thank you very much for this comment. We rewrote the introduction, including these
796 lines. In those lines of the original manuscript, we regarded the restoring advection feedback as a
797 fast process and linked it to period therein, although we also discussed its destabilizing effect
798 elsewhere in the manuscript. Actually, throughout the paper, we separate the discussion of
799 temperature's effects on period from that on stability.

800

801 *11. line 86: change the word "unravel".* The introduction has been rewritten.

802 *12. line 106: define abbreviation before use.* This abbreviation is defined in the revised
803 manuscript.

804 *13. formulation of the box models: it is assumed that deep water completely upwells in the interior*
805 *of the ocean. However, people have also found that upwelling in the Southern Ocean is also, if*
806 *not more, important than the interior diffusive upwelling. How important the authors think the*
807 *Southern Ocean is for their results here?*

808 **Responses:** Thank you very much for this comment. In the rewritten introduction, the Southern
809 Ocean and even the Arctic Ocean may play a role in the multicentennial AMOC oscillation.
810 Therefore, these two oceans definitely play a part. The box models can only grab the most
811 fundamental physics, the importance of the Southern Ocean cannot be specified in the one-
812 hemisphere model.

813 Actually, we are working on a two-hemisphere model, in which the Southern Ocean is
814 included. We believe the Southern Ocean plays a role in the multicentennial oscillation of the
815 AMOC, and also in the millennial oscillation of the climate system. We are also using a coupled
816 climate model to study the role of the Southern Ocean.

817 Our guess is that including the Southern Ocean in the box model should prolong the oscillation
818 period of the AMOC. The system stability might also be affected, since q' now is likely to be
819 determined by the density difference between the North Atlantic and the Southern Ocean. However,
820 the temperature and salinity advection feedbacks related to the original box 2, together with the
821 northern subpolar vertical mixing process, will still be at work; thus, the basic self-sustained
822 oscillation mechanism should remain the same.

823

824 14. Usually τ is a symbol for timescale. Here it is a reversal of a timescale. Using another letter
825 may be less confusing to some people. Revised. τ is replaced with γ throughout.

826 15. Fw: directly say that it is an artificial salt flux representing the effect of freshwater flux, in
827 context and in the table. Revised.

828 16. the choice of 10 Sv for q is not very characteristic of time-averaged AMOC. Why is this value
829 chosen. In the discussion part of the paper, we see some references in which the timescale of
830 MOC oscillation is dependent on q . How different the results of this manuscript will be if q is
831 more realistic? Also, q is directly given, is it consistent with the time-mean large-scale
832 meridional density gradient as specified by the formula (related by λ)?

833 **Responses:** Thank you very much for this comment.

834 We admit that the choice of \bar{q} is critical to the oscillation timescale. In the observation, the
835 maximum mean AMOC is about 20 Sv. This mass transport includes water in the upper 1000 m and
836 water from the Southern Hemisphere. In a one-hemispheric box model with the upper-ocean depth
837 set to 500 m, the mean mass transport should be remarkably smaller than the realistic value.
838 Otherwise, the turnover timescale for a one-hemisphere box model would be unrealistically short. \bar{q}
839 in GT95 is much larger than ours; and their subpolar boxes are much smaller than ours, so that the
840 dominant timescale in GT95 box model is the decadal timescale, instead of the centennial timescale.
841 Our single-hemispheric model incorporates only the AMOC recirculating in the Northern
842 Hemisphere, so that a smaller mean AMOC is reasonable.

843 Table R1 lists the eigenvalues for the 4TS and 4S models under the parameters in Table 1, but
844 with $\bar{q} = 18 Sv$. The e-folding time is 36 years for decaying oscillations in the 4TS model, while it
845 is 40 years for decaying oscillations in the 4S model, both are more stable than their counterparts
846 under $\bar{q} = 10 Sv$. The period is 208 years in the 4TS model, while it is 220 years in the 4S model,
847 both shorter than their counterparts under $\bar{q} = 10 Sv$. The results suggest that a larger \bar{q} leads to a
848 more stable system with a faster oscillation, consistent with the results in section 5b of this revised
849 manuscript. However, the overall results of this paper are not substantially altered with a larger \bar{q}
850 and there is still the potential for multicentennial oscillation.

851 TABLE R1. Eigenvalues ($10^{-10} s^{-1}$) for the 4TS and 4S models under the parameters of Table 1,
852 but with $\bar{q} = 18 Sv$.

4TS	4S	Physical Significance
$-8.79 \pm 9.57i$	$-7.91 \pm 9.07i$	Oscillatory mode

0	0	Zero mode
-393	—	Damped mode
-330	—	Damped mode
-65.5	-65.5	Damped mode
-9.84	—	Damped mode
-1.47	—	Damped mode

853

854 In our parameterization, it is q' rather than \bar{q} that is proportional to the meridional difference
855 of density anomaly $\Delta\rho'$. Therefore, \bar{q} is directly given instead of being related to the large-scale
856 meridional difference of the mean density.

857

858 *17. The upper layer where restoring happens is as thick as 500 meters, which is too deep for*
859 *restoring to be justifiable. Also, the total thickness of ocean in these box models is 4 km, but*
860 *these models are only for the upper cell which may not penetrate so deep in the ocean.*

861 **Responses:** Thank you very much for this comment. We admit that in reality the restoring can only
862 happen within the surface layer with a depth of a few tens of meters, which is also seen in studies
863 using models of higher complexity (Marotzke and Willebrand 1991; Weaver and Sarachik 1991;
864 Mysak et al. 1993). However, there is not such surface layer in our theoretical model due to its
865 simplicity.

866 For model conciseness, we permit the restoring to happen within the whole depth range of
867 boxes 1 and 2, which is commonly adopted in box model studies. As compensation, we chose 1-
868 year restoring timescale, which is significantly long due to the rather thick restoring surface layer
869 (D_1 , 500 m). In GT95, the authors permitted the surface restoring to take place over the 300-m thick
870 surface layer, and chose 180-day restoring timescale. Also, in RT97 the surface restoring can
871 happen over the 1000-m thick surface layer; they also chose 1-year restoring timescale to
872 compensate such thickness.

873 Boxes 1 and 2 are for the upper ocean, while boxes 3 and 4 are for the deeper ocean including
874 the southward NADW, not the upper cell. This is also evidenced by the southward mean AMOC
875 between boxes 3 and 4 in our theoretical model. Therefore, the 4-km depth for our box model is
876 reasonable.

877

878 **References:**

879 Griffies, S. M., and E. Tziperman, 1995: A linear thermohaline oscillator driven by stochastic atmospheric
880 forcing. *J. Climate*, 8, 2440-2453.

881 Marotzke, J., and J. Willebrand, 1991: Multiple equilibria of the global thermohaline circulation. *J. Phys.*
882 *Oceanogr.*, 21, 1372-1385.

883 Mysak, L. A., T. F. Stocker, and F. Huang, 1993: Century-scale variability in a randomly forced, two-
884 dimensional thermohaline ocean circulation model. *Climate Dyn.*, 8, 103-116.

885 Rivin, I., and E. Tziperman, 1997: Linear versus self-sustained interdecadal thermohaline variability in a
886 coupled box model. *J. Phys. Oceanogr.*, 27, 1216-1232.

887 Weaver, A. J., and E. S. Sarachik, 1991: Evidence for decadal variability in an ocean general-circulation
888 model - An advective mechanism. *Atmos. Ocean*, 29, 197-231.

889

890 *18. How are the parameters T_1 , T_2 , S_1 , S_2 , T_1^* and T_2^* determined? Are they based on some*
891 *dataset?*

892 **Responses:** Thank you very much for this comment. They are based on the CESM1 simulation
893 analyzed in LY22.

894

895 *19. equation 3: The use of letter κ_0 may confuse it with diffusivity. What is the unit of this*
896 *coefficient, and how is it determined?*

897 **Responses:** Thank you very much for this comment. In the revised manuscript, we rewrote this
898 paragraph.

899 κ_0 denotes the restoring coefficient, which reflects the strength of SST relaxation toward T_1^* or
900 T_2^* , with the units of $W/(m^2 \cdot ^\circ C)$. Originally, it was determined from observations. Bretherton
901 (1982) suggested that a large value like $100 W/(m^2 \cdot ^\circ C)$ be chosen for the sea surface with a few
902 tens of kilometers across, since atmospheric advection of heat is also important for small-scale
903 surface heat flux. He also stated that a small value like $2 W/(m^2 \cdot ^\circ C)$ be chosen for global SST,
904 whose restoring is largely determined by local radiation instead of atmospheric advection. In later
905 modelling studies, this value was usually set to a constant, leading to 1-2 month restoring timescale
906 for a surface layer of a few tens of meters (Marotzke and Willebrand 1991; Weaver and Sarachik
907 1991; Mysak et al. 1993; Pierce 1996).

908

909 **References:**

910 Bretherton, F. P., 1982: Ocean climate modeling. *Prog. Oceanogr.*, 11, 93-129.

911 Marotzke, J., and J. Willebrand, 1991: Multiple equilibria of the global thermohaline circulation. *J. Phys. Oceanogr.*, 21, 1372-1385.

913 Mysak, L. A., T. F. Stocker, and F. Huang, 1993: Century-scale variability in a randomly forced, two-dimensional thermohaline ocean circulation model. *Climate Dyn.*, 8, 103-116.

915 Pierce, D. W., 1996: Reducing phase and amplitude errors in restoring boundary conditions. *J. Phys. Oceanogr.*, 26, 1552-1560.

917 Weaver, A. J., and E. S. Sarachik, 1991: Evidence for decadal variability in an ocean general-circulation model - An advective mechanism. *Atmos. Ocean*, 29, 197-231.

919

920 *20. line 157: for equation of state, how realistic is it to use a uniform value for α and β across such a wide latitudinal range? Around what temperature and salinity are these coefficients determined?*

923 **Responses:** Thank you very much for this comment.

924 The thermal expansion coefficient of $1.468 \times 10^{-4} \text{ }^\circ\text{C}^{-1}$ and the haline contraction efficient
925 of $7.61 \times 10^{-4} \text{ } \text{psu}^{-1}$ are derived from UNESCO (1987) around 9°C and $35 \text{ } \text{psu}$ under 0 dbar. We
926 admit that these two coefficients can vary not only with temperature, but also with salinity and
927 pressure. However, we chose constant values in our study for simplicity, as the essence for a box
928 model lies in its ability of revealing mechanisms via its conciseness instead of statistical accuracy.
929 Constant thermal expansion and haline contraction coefficients have been adopted in a large
930 number of theoretical models for simplicity, and some even covered a larger latitudinal extent
931 (Scott et al. 1999; Lucarini and Stone 2005; Alkhayuon et al. 2019; Shi and Yang 2021; Wei and
932 Zhang 2022).

933

934 **References:**

935 Alkhayuon, H., P. Ashwin, L. C. Jackson, C. Quinn, and R. A. Wood, 2019: Basin bifurcations, oscillatory
936 instability and rate-induced thresholds for Atlantic Meridional Overturning Circulation in a global
937 oceanic box model. *Proc Math Phys Eng Sci*, 475, 20190051.

938 Lucarini, V., and P. H. Stone, 2005: Thermohaline circulation stability: A box model study. Part I:
939 Uncoupled model. *J. Climate*, 18, 501-513.

940 Scott, J. R., J. Marotzke, and P. H. Stone, 1999: Interhemispheric thermohaline circulation in a coupled box
941 model. *J. Phys. Oceanogr.*, 29, 351-365.

942 Shi, J. Q., and H. J. Yang, 2021: Bjerknes compensation in a coupled global box model. *Climate Dyn.*, 57,
943 3569-3582.

944 UNESCO, 1987. International oceanographic tables. Tech. Paper Mar., Sci. 40, 196 pp.

945 Wei, X., and R. Zhang, 2022: A simple conceptual model for the self-sustained multidecadal AMOC
946 variability. *Geophys. Res. Lett.*, 49, 11.

947

948 *21. line 196–197: “1025 years for growing oscillations” and “576 years for decaying oscillations”*
949 *instead of “positive 1025 years” and “negative 576 years”. Revised.*

950 *22. line 199: Figure 3 is mentioned before Figure 2. Revised.*

951 *23. line 208: change “ $y=0$ ” to “ $Im(\omega) = 0$ ” to be more clear. Revised.*

952 *24. Figure 2(a): The changes between the “S-only” and “T” models happen mainly where λ is*
953 *large. but less so when it is small. Why?*

954 **Responses:** Thank you very much for this comment. Compared to salinity-only models, the range
955 of λ with complex eigenvalues is increased in the TS models, especially when λ is large. As λ
956 increases, the imaginary parts of the S models will move to 0 sooner than those in the TS models;
957 thus, the changes between the S and TS models are the most evident when λ is larger.

958

959 *25. line 225 (Figure 2 caption): “zoomed-in” instead of “magnified”. Revised.*

960 *26. table 3: Why would λ be a negative value? Will it actually be zero? If so, what does it*
961 *physically imply?*

962 **Responses:** Thank you very much for this comment. Physically, λ cannot be zero or negative since
963 the AMOC will be stronger under a bigger meridional difference of density anomaly (positive λ).
964 However, we decided to include $\lambda \leq 0$ in Fig. 2 so as to show both the upper and lower limits for
965 the existence of imaginary part. This would not influence the results presented in this manuscript
966 since the analyses are only related to $\lambda > 0$ case.

967

968 *27. line 248: But the advection feedback is related to large-scale circulation, why would it be local?*

969 **Responses:** Thank you very much for this comment. The local advection feedback we referred to
970 here is the negative temperature advection feedback $q'(\overline{T}_1 - \overline{T}_2)$. It influences T_2 without affecting
971 T_1 ; therefore, it is a local one. Also in reality and complex coupled model, q' denotes local change.
972 When it comes to the feedback between mean advection of temperature anomalies and AMOC
973 anomaly ($\overline{q}(T'_1 - T'_2)$), it influences T_2 through transporting temperature anomaly from box 1 to box
974 2; thus, it is a remote feedback (non-local).

975

976 *28. line 272–277: Move this paragraph to before the discussion of Figure 2. Revised.*

977 *29. line 293–294: Temperature is still advected and enters the EOS, correct? I am a little confused*
978 *why the $\tau = 0$ case should be the same as the salinity-only case.*

979 **Responses:** Thank you very much for this comment. Temperature will always be advected as long
980 as the temperature equations are involved. However, when $\gamma = 0$ (or approaches 0, more
981 specifically), from Eq. (8) we have $\overline{T}_1 - \overline{T}_2 = 0$, thus $\overline{T}_1 = \overline{T}_2 = \overline{T}_3 = \overline{T}_4$. Therefore, $q'(\overline{T}_1 - \overline{T}_2)$,
982 $q'(\overline{T}_2 - \overline{T}_3)$, $q'(\overline{T}_3 - \overline{T}_4)$, and $q'(\overline{T}_4 - \overline{T}_1)$ will be 0; thus, there will be no variation in temperature
983 at all. Consequently, the temperature advection would not warm or cool any boxes; and the
984 temperature equations are useless.

985

986 *30. line 306: What is the reasonable range of τ since this is an artificial form of forcing? How is*
987 *this reasonable range determined?*

988 **Responses:** Thank you very much for this comment. Typically, the restoring timescale is 1-2
989 months with a surface layer of a few tens of meters. However, since our theoretical model does not
990 have such an explicit surface layer, we permit the restoring to happen over the whole depth range of
991 boxes 1 and 2. Consequently, the surface layer now is 500 m, almost 10 times the typical value. As
992 compensation, we should choose a restoring timescale much longer than the typical value of 1-2
993 months. A value of 1-year or longer is reasonable, while 9-10 months might also work. However, it
994 should not be as short as a season or shorter.

995

996 *31. lines 359–361: Judging from Figure 4, over some range of τ 3TS is more stable. Be more*
997 *specific in this sentence.*

998 **Responses:** Thank you very much for this comment. As we have explained earlier in our replies, the
999 reasonable range of restoring timescale is no shorter than one year and also cannot be too long.

1000 Therefore, although over some range of γ the 3TS model is more stable, the 4TS model is more
1001 stable than the 3TS one under realistic γ .

1002

1003 *32. line 394–396: Why? It is the averaged T/S in box 2 and 3 that determines the strength of MOC,*
1004 *which will not be changed here since mixing only occurs between these two boxes. Is this*
1005 *stabilizing effect from the surface conditions?*

1006 **Responses:** Thank you very much for this comment. A warmer (more saline) anomaly in the
1007 subpolar region has stabilizing (destabilizing) effect on AMOC oscillation. The subpolar vertical
1008 mixing moves the anomalies from box 2 to box 3 faster. Therefore, although the mixing process
1009 conserves salt and temperature, it removes the warm and salty anomalies from the subpolar region
1010 more rapidly, limiting their stabilizing and destabilizing effects, respectively.

1011

1012 *33. section 4b: The only purpose of this is to make self-sustained oscillation possible by using some*
1013 *complicated relation between density gradient and MOC, otherwise I find it distracting.*
1014 *Consider putting in supplementary.*

1015 **Responses:** Thank you very much for this comment.

1016 We have improved the logic of section 4 in the revised manuscript. Section 4a is to show that
1017 the self-sustained oscillation is able to be realized even with the presence of destabilizing vertical
1018 mixing (the destabilizing temperature mixing overcomes the stabilizing salinity mixing).

1019 In section 4b, we use this seemingly complicated relation between meridional difference of
1020 density anomaly and AMOC anomaly (in fact this relation is physically simple; it means restraining
1021 of q' is introduced when the meridional difference of density anomaly is large) to introduce a
1022 degree of nonlinearity in the absence of mixing. The realization of self-sustained oscillation in both
1023 the 4TS and 3TS models here highlights the role of nonlinearity, while the self-sustained oscillation
1024 is not sensitive to the exact form of the restraining term. Only with what in section 4b, can we
1025 generalize the LY22 self-sustained oscillation mechanism of “a combination of salinity advection
1026 and enhanced mixing” to “a linearly growing oscillation dominated by advection and a nonlinear
1027 restraining effect from restraining terms” in this manuscript. Therefore, we prefer to place section
1028 4b in the main body of the manuscript.

1029

1030 *34. line 455: “much less stable than ours” instead of “stability far lower than ours”.* This
1031 paragraph is deleted.

1032 35. *line 563: What do you mean by “random components”?* We deleted discussion on this.

1033 36. *line 568: “models of higher complexity” instead of “higher complexity models”.* This sentence
1034 is deleted.

1035 37. *line 576–579: The collapse is more related to bifurcation than stability of oscillation.*

1036 **Responses:** Thank you very much for this suggestion. We have removed the contents about AMOC
1037 collapse.

1038



Click here to access/download

Additional Material for Reviewer Reference

THC_MultiCentennial_Theory_P2_tracked_20230130.pdf

f