

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

3
4 Kunpeng Yang^{1,2}, Haijun Yang*^{1,2} and Yang Li³

5
6 ¹*Department of Atmospheric and Oceanic Sciences and Institute of Atmospheric Science and CMA-*
7 *FDU Joint Laboratory of Marine Meteorology, Fudan University, Shanghai, 200438, China.*

8 ²*Shanghai Scientific Frontier Base for Ocean-Atmosphere Interaction Studies, Fudan University,*
9 *Shanghai 200438, China.*

10 ³*Department of Atmospheric and Oceanic Sciences, School of Physics, Peking University, Beijing,*
11 *100871, China.*

12
13
14
15
16 *Journal of Climate*

17 Submitted

18 October 5, 2022

19
20 *Corresponding author address: Haijun Yang, Department of Atmospheric and Oceanic Sciences,
21 Fudan University, 2005 Songhu Road, Shanghai China, 200438.

22 Email: yanghj@fudan.edu.cn.

23

24

Abstract

25 In the first part of our research on self-sustained multicentennial oscillation of the Atlantic
26 meridional overturning circulation (AMOC), we proposed a hemispheric box model considering only
27 the saline process. In this paper, we consider both thermal and saline processes in the box model and
28 employ mixed boundary conditions, so as to include more realistic physics. Generally, the thermal
29 process has a stabilizing effect on the system, and additional physics, such as enhanced subpolar
30 mixing or a nonlinear relation between the AMOC and meridional density gradient, is still needed to
31 realize a self-sustained oscillation. Specifically, the thermal process exerts mainly two effects on the
32 system: shortening marginally the multicentennial oscillation period of the AMOC, and stabilizing the
33 oscillating system and subpolar stratification, which are contributed by the fast surface temperature
34 restoring, the negative temperature-advection feedback and subpolar temperature stratification,
35 respectively. The oscillation properties are controlled by the balance of destabilizing salinity
36 advection and stabilizing temperature advection. Different from salinity-only situation, the enhanced
37 subpolar mixing in the current situation makes the system more unstable. Weaker meridional
38 temperature gradient and stronger meridional salinity gradient can lead to weaker temperature-
39 advection feedback and stronger salt-advection feedback, and thus a longer AMOC oscillation period
40 with less stability at multicentennial timescale, which might be expected in the future due to more
41 intense high-latitude warming and freshwater hosing.

42 **Keywords:** AMOC, Self-sustained oscillation, Temperature feedbacks, Box model

43

44 **1. Introduction**

45 In our first publication on the multicentennial oscillation of the Atlantic meridional overturning
46 circulation (AMOC) (Li and Yang 2022, hereafter LY22), we used a 4-box model to study the self-
47 sustained multicentennial AMOC oscillation, which only considered salinity equations. Here, we
48 expand our horizon to explore theoretically the multicentennial oscillation of the AMOC by including
49 both thermal and saline processes. Thus, the term AMOC in this paper refers to its thermohaline
50 circulation portion that is controlled by both temperature and salinity. It is well recognized that the
51 AMOC is a crucial regulator for the North Atlantic and likely the global climate over a wide range of
52 timescales (Chabaud et al. 2014; Muir and Fedorov 2015; Zhang et al. 2019). Low-order theoretical
53 models are essentials for understanding the AMOC dynamics. To isolate the most fundamental
54 dynamics, only salinity equations were kept in the theoretical model of LY22, paralleling a few other
55 studies (Rahmstorf 1996; Cimatoribus et al. 2014; Sévellec and Fedorov 2014). However, such
56 treatment is unphysical to an extent since all the temperature-related effects were excluded.

57 In ocean-only theoretical models of the AMOC, mixed boundary conditions (Haney 1971) are
58 often employed; that is, sea surface temperature (SST) is restored to a prescribed value following the
59 Newtonian law, while sea surface salinity (SSS) is forced by the surface freshwater flux. As seen in
60 quite a few studies, the negative temperature-advection feedback (Stommel 1961; Walin 1985) and
61 the positive restoring-advection feedback (Griffies and Tziperman 1995, hereafter GT95; Scott et al.
62 1999; Colin de Verdière 2010) are included in their temperature equations. The former works as
63 follows: commencing with an initial positive AMOC perturbation, the elevated poleward heat
64 transport reduces subpolar density, therefore restraining deep-water formation, which is followed by
65 the slowdown of the AMOC. The latter works as an opponent to the former: an enhanced AMOC
66 increases the subpolar SST, leading to restore the warming itself, and the reduced positive SST
67 anomaly hinders the AMOC slowdown. The restoring-advection feedback strengthens as the restoring
68 process (timescale) speeds up (shortens). However, this feedback merely offsets, but never overruns
69 the temperature-advection feedback even under extremely strong restoring. Hence, the net effect of
70 temperature feedbacks is to stabilize the system. It has also been illustrated in other studies that the
71 oceanic thermal process influences the stability of AMOC system (Zhang et al. 1993; Marotzke and
72 Stone 1995; Rahmstorf and Willebrand 1995; Marotzke 1996). Consequently, adding temperature
73 equations to the salinity-only model of LY22 should be more realistic for the multicentennial AMOC
74 oscillation.

75 Usually, the thermal process is faster than the saline process due to the fast SST restoring, whose
76 timescale is not more than a few years (Pierce 1996). It is thus intuitive that including this fast process
77 might shorten the multicentennial period. However, Schmidt and Mysak (1996) considered that such
78 fast restoring removes high-frequency anomalies and therefore might prolong the multicentennial
79 period. They then stated that such lengthening is not so obvious; and their main focus was on the
80 system's stability, leaving this question unanswered. Therefore, it would be intriguing to see the role
81 of temperature in influencing AMOC oscillation timescale. An in-depth understanding of temperature
82 effects in shaping multicentennial AMOC oscillation should also provide us insight about how the
83 AMOC will respond to future climate change.

84 In this study, we extend the 4-box salinity-only model in LY22 to a temperature-salinity one; and
85 we employ mixed boundary conditions. The conciseness of this model enables stability analysis. We
86 aim to unravel the effects induced by thermal process. We then work on realization of a self-sustained
87 oscillation with bounding terms affiliated, in order to test whether the advection process dominates
88 the eigenmode. Motivated by the difference in period and stability between our model and other
89 studies, we also examine the sensitivity of eigenmode to model parameters controlling the model
90 geometry, flow properties and feedback processes.

91 This paper is structured as follows. In section 2, a 4-box temperature-salinity model (hereafter
92 4TS) is introduced, followed by illustration of temperature and salinity feedbacks. In section 3, the
93 role of temperature equations is analyzed. In section 4, we test two ways for realizing a self-sustained
94 oscillation, and elucidate the critical role of advection process. In section 5, the sensitivities of
95 eigenmode's period and stability to model parameters are examined. Summary and discussion are
96 presented in section 6.

97

98 **2. Box model**

99 *a. Model formulae*

100 The model we use here is a hemispheric 4-box model with identical geometry to that in LY22
101 (Fig. 1a). The model domain is 60° in longitude, with the tropical and subpolar boxes spanning over
102 0° - 45° N and 45° - 70° N, respectively. The AMOC flows through the boxes in a clockwise sense.
103 Excluding multi-equilibria, we do not discuss possibility of a reversed AMOC cell. Analogous box
104 models have been widely used (Joyce 1991; Huang et al. 1992; GT95). In the 4-box salinity-only

105 model (hereafter 4S) of LY22, only salinity equations were used to obtain analytical solutions. In the
 106 4TS model, we have:

$$107 \quad V_1 \dot{T}_1 = q(T_4 - T_1) + V_1 \tau(T_1^* - T_1) \quad (1a)$$

$$108 \quad V_2 \dot{T}_2 = q(T_1 - T_2) + V_2 \tau(T_2^* - T_2) \quad (1b)$$

$$109 \quad V_3 \dot{T}_3 = q(T_2 - T_3) \quad (1c)$$

$$110 \quad V_4 \dot{T}_4 = q(T_3 - T_4) \quad (1d)$$

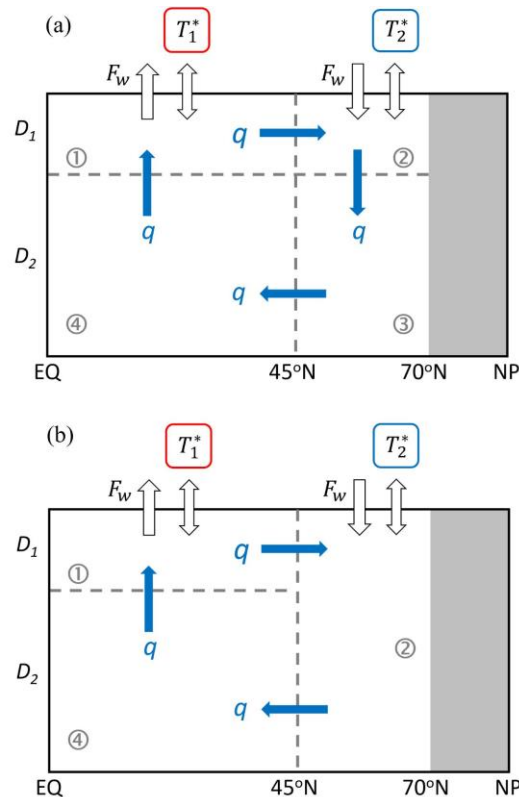
$$111 \quad V_1 \dot{S}_1 = q(S_4 - S_1) + F_w \quad (1e)$$

$$112 \quad V_2 \dot{S}_2 = q(S_1 - S_2) - F_w \quad (1f)$$

$$113 \quad V_3 \dot{S}_3 = q(S_2 - S_3) \quad (1g)$$

$$114 \quad V_4 \dot{S}_4 = q(S_3 - S_4) \quad (1h)$$

115 It is an advection-dominated box model, with mixed boundary conditions where Haney-style
 116 restoring for SST (Haney 1971) and surface freshwater flux for SSS are adopted. V_i , T_i and S_i are the
 117 volume, temperature and salinity in each box. Dots over variables denote their temporal derivatives. q
 118 stands for AMOC strength. F_w represents the surface freshwater flux, which actually takes the form of
 119 virtual salt flux. T_1^* and T_2^* correspond to the restoring temperatures for boxes 1 and 2, respectively.
 120 The Newtonian restoring coefficient τ is also the reciprocal of the restoring timescale for T_1 and T_2 .



121

122 FIG. 1. Schematic diagrams of temperature-salinity box models. (a) The 4-box model. (b) The 3-box model
 123 reduced from the 4-box one. The circled numbers ①, ②, ③, and ④ denote the ocean boxes. Boxes 1 and 4
 124 stand for the upper and lower tropical oceans, respectively, while boxes 2 and 3 stand for the upper and lower
 125 subpolar oceans, respectively. D_1 and D_2 represent the upper and lower ocean depths, respectively. The net
 126 freshwater flux out of (into) the tropical (subpolar) ocean is represented by F_w . T_1^* (T_2^*) is the restoring
 127 temperature of the tropical (subpolar) ocean. q represents the AMOC.
 128

129 The equilibrium solutions at steady state can be easily written as:

$$130 \quad \overline{T}_1 = T_1^* - \frac{\overline{q}V_2(T_1^* - T_2^*)}{\overline{q}(V_1 + V_2) + V_1V_2\tau}, \quad \overline{T}_2 = \frac{V_1T_1^* + V_2T_2^* - V_1\overline{T}_1}{V_2} = \overline{T}_3 = \overline{T}_4 \quad (2a)$$

$$131 \quad \overline{S}_1 = F_w/\overline{q} + \overline{S}_2, \quad \overline{S}_2 = \overline{S}_3 = \overline{S}_4 \quad (2b)$$

132 Variables with overbar denote their equilibrium values. Similar to previous studies (Winton and
 133 Sarachik 1993; Cessi 1994; Roebber 1995), we choose $\overline{q} = 10 Sv$. Following LY22, the upper and
 134 total ocean depths are 500 and 4000 m, respectively. We set \overline{S}_2 and F_w to 33.5 psu and 25.0 psu · Sv,
 135 respectively, leading to $\overline{S}_1 = 36 psu$. Straightway, T_1^* and T_2^* are considered to be close to the
 136 averaged realistic SSTs in the tropical and subpolar regions; thus we have 25°C and 7°C, respectively.
 137 The setting of τ uses the idea of Bretherton (1982), and the restoring timescale $1/\tau$ is represented as
 138 follows:

$$139 \quad 1/\tau = \frac{\rho_w c \Delta z A}{\kappa_0 A} = \frac{\rho_w c \Delta z}{\kappa_0} \quad (3)$$

140 Here, A is the area of ocean surface. ρ_w and c are the typical density and specific heat of seawater, set
 141 to 1027 kg · m⁻³ and 3850 J/(kg · °C), respectively. A value of 30 m is given to the thickness of the
 142 surface layer Δz (not the upper ocean). κ_0 is a restoring coefficient. Bretherton (1982) stated that κ_0
 143 should be small when averaged over the whole globe, while it is larger if a smaller area is considered.
 144 In view of our one-hemisphere configuration, it is reasonable to choose $\kappa_0 = 4 W/(m^2 · °C)$, yielding
 145 a restoring timescale of almost one year in our box model. For convenience, we set $\tau =$
 146 $3.171 \times 10^{-8} s^{-1}$, corresponding to a precise 1-year restoring timescale.

147 The total AMOC strength q could be separated into an equilibrium portion \overline{q} and an anomalous
 148 portion q' . We consider a linear relation between q' and thickness-weighted meridional density
 149 gradient anomaly $\Delta\rho'$. Both q' and $\Delta\rho'$ can be decomposed into temperature-driven portion ($q'_T, \Delta\rho'_T$)
 150 and salinity-driven portion ($q'_S, \Delta\rho'_S$); therefore, we have:

$$151 \quad q = \bar{q} + q' \quad (4a)$$

$$152 \quad q' = q'_T + q'_S = \lambda \Delta \rho'_T + \lambda \Delta \rho'_S = \lambda \Delta \rho' \quad (4b)$$

153 where

$$154 \quad \Delta \rho'_T = -\rho_0 \alpha [\delta (T'_2 - T'_1) + (1 - \delta)(T'_3 - T'_4)] \quad (4c)$$

$$155 \quad \Delta \rho'_S = \rho_0 \beta [\delta (S'_2 - S'_1) + (1 - \delta)(S'_3 - S'_4)] \quad (4d)$$

$$156 \quad \delta = \frac{V_1}{V_1 + V_4} = \frac{V_2}{V_2 + V_3} = \frac{D_1}{D} \quad (4e)$$

157 The sensitivity of q' to $\Delta \rho'$ is represented by a linear closure coefficient λ . ρ_0 , α and β are the
 158 reference density, thermal expansion and haline contraction coefficients for seawater, respectively. D_1
 159 and D correspond to the upper and total ocean depths, respectively. T'_i and S'_i are the temperature and
 160 salinity anomalies of box i , respectively. A summary of the standard parameter values is provided in
 161 Table 1.

162 TABLE 1. Standard values of the parameters used.

Symbol	Physical Significance	Value
V_2	Volume of box 2	$2.8 \times 10^{15} \text{ m}^3$
V_1, V_3, V_4	Volumes of boxes 1, 3 and 4, respectively	$5V_2, 7V_2, 35V_2$
D_1, D_2, D	Thicknesses of the upper, lower oceans and the entire ocean	500, 3500, 4000 m
T_1^*, T_2^*	Restoring temperatures of boxes 1 and 2	25°C, 7°C
τ	Restoring coefficient of boxes 1 and 2	$3.171 \times 10^{-8} \text{ s}^{-1}$
$\bar{S}_1, \bar{S}_2, \bar{S}_3, \bar{S}_4$	Equilibrium salinities of boxes 1, 2, 3 and 4	36, 33.5, 33.5, 33.5 psu
\bar{q}	Equilibrium strength of AMOC	10 Sv ($10^6 \text{ m}^3 \text{ s}^{-1}$)
F_w	Surface freshwater flux	25.0 $psu \cdot Sv$
λ	Linear closure coefficient	12 $Sv \cdot kg^{-1} \text{ m}^3$
ρ_0	Reference seawater density	$1.00 \times 10^3 \text{ kg m}^{-3}$
α	Thermal expansion coefficient	$1.468 \times 10^{-4} \text{ }^\circ\text{C}^{-1}$
β	Haline contraction coefficient	$7.61 \times 10^{-4} \text{ psu}^{-1}$

164 We linearize Eq. (1) as follows:

$$165 \quad V_1 \dot{T}'_1 = q'(\overline{T}_4 - \overline{T}_1) + \overline{q}(T'_4 - T'_1) - V_1 \tau T'_1 \quad (5a)$$

$$166 \quad V_2 \dot{T}'_2 = q'(\overline{T}_1 - \overline{T}_2) + \overline{q}(T'_1 - T'_2) - V_2 \tau T'_2 \quad (5b)$$

$$167 \quad V_3 \dot{T}'_3 = \overline{q}(T'_2 - T'_3) \quad (5c)$$

$$168 \quad V_4 \dot{T}'_4 = \overline{q}(T'_3 - T'_4) \quad (5d)$$

$$169 \quad V_1 \dot{S}'_1 = q'(\overline{S}_4 - \overline{S}_1) + \overline{q}(S'_4 - S'_1) \quad (5e)$$

$$170 \quad V_2 \dot{S}'_2 = q'(\overline{S}_1 - \overline{S}_2) + \overline{q}(S'_1 - S'_2) \quad (5f)$$

$$171 \quad V_3 \dot{S}'_3 = \overline{q}(S'_2 - S'_3) \quad (5g)$$

$$172 \quad V_4 \dot{S}'_4 = \overline{q}(S'_3 - S'_4) \quad (5h)$$

173 In LY22, we assumed an extremely strong vertical mixing between subpolar boxes 2 and 3; the
 174 4S model can be reduced to a 3-box salinity-only model (hereafter 3S). Applying the same treatment
 175 to the 4TS model, a 3-box temperature-salinity model (hereafter 3TS; Fig. 1b) can be obtained. Now,
 176 Eqs. (4c-e) become,

$$177 \quad \Delta \rho'_T = -\rho_0 \alpha [T'_2 - \delta T'_1 - (1 - \delta) T'_4] \quad (6a)$$

$$178 \quad \Delta \rho'_S = \rho_0 \beta [S'_2 - \delta S'_1 - (1 - \delta) S'_4] \quad (6b)$$

$$179 \quad \delta = \frac{V_1}{V_1 + V_4} = \frac{D_1}{D} \quad (6c)$$

180 and Eqs. (5a-h) are reduced to:

$$181 \quad V_1 \dot{T}'_1 = q'(\overline{T}_4 - \overline{T}_1) + \overline{q}(T'_4 - T'_1) - V_1 \tau T'_1 \quad (7a)$$

$$182 \quad V_2 \dot{T}'_2 = q'(\overline{T}_1 - \overline{T}_2) + \overline{q}(T'_1 - T'_2) - V_2 \tau T'_2 \quad (7b)$$

$$183 \quad V_4 \dot{T}'_4 = \overline{q}(T'_2 - T'_4) \quad (7c)$$

$$184 \quad V_1 \dot{S}'_1 = q'(\overline{S}_4 - \overline{S}_1) + \overline{q}(S'_4 - S'_1) \quad (7d)$$

$$185 \quad V_2 \dot{S}'_2 = q'(\overline{S}_1 - \overline{S}_2) + \overline{q}(S'_1 - S'_2) \quad (7e)$$

$$186 \quad V_4 \dot{S}'_4 = \overline{q}(S'_2 - S'_4) \quad (7f)$$

187

188 *b. Stability analysis*

189 Let us first examine the eigenvalues of the 4TS model. Table 2 lists the eight eigenvalues of Eq.
 190 (5) using the parameters in Table 1. The eigenvalues in the 4S model of LY22 using the same
 191 parameters are listed in Table 2 for comparison.

192 TABLE 2. Eigenvalues (10^{-10} s^{-1}) for the 4TS and 4S models using the parameters of Table 1.

4TS	4S	Physical significance
$-0.55 \pm 6.59i$	$0.31 \pm 5.83i$	Oscillatory mode
0	0	Zero mode
-366	—	Damped mode
-324	—	Damped mode
-37.4	-37.4	Damped mode
-5.28	—	Damped mode
-0.78	—	Damped mode

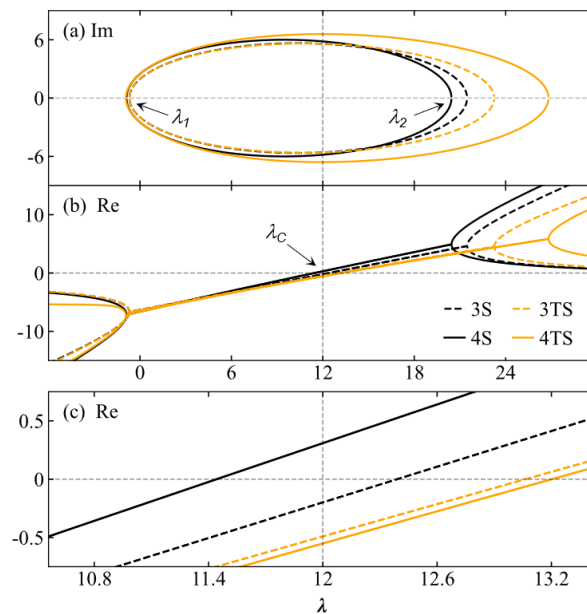
193

194 There is still a pair of conjugate eigenvalues ($-0.55 \pm 6.59i$) in the 4TS model. The weakly
 195 unstable oscillatory mode ($0.31 \pm 5.83i$) in the 4S model becomes a weakly damped oscillatory mode
 196 in the 4TS model (Fig. 3a). That is, the e-folding time changes from positive 1025 years in the 4S
 197 model to negative 576 years in the 4TS model, and the period changes from 340 years to about 300
 198 years. This seems to suggest that the thermal processes have a stabilizing effect on the system, and
 199 tend to shorten the oscillation period slightly (Fig. 3a). The zero mode (eigenvalue 0) represents the
 200 climatological mean state. The other five eigenvalues in the 4TS model represent five purely damped
 201 modes, which are not of our concern.

202 The stability of the box model system is strongly dependent on the linear closure parameter λ ,
 203 i.e., the sensitivity of the AMOC to the meridional density gradient as formulated in Eq. (4b). The
 204 critical role of λ and its physical explanation can be found in LY22. In this paper, we simply solve
 205 Eqs. (5) and (7) numerically to investigate how λ affects the stabilities of the 4TS and 3TS models.
 206 Figure 2 shows the dependence of real and imaginary parts of the oscillatory mode on λ . The results
 207 from the 4S and 3S models in LY22 are also plotted in Fig. 2 for comparison. The intersections

208 between line $y = 0$ and the stability diagrams of each model, $(\lambda_c, 0)$, $(\lambda_1, 0)$ and $(\lambda_2, 0)$, correspond
 209 to the instability threshold (Fig. 2b), the lower and upper limits for the existence of the imaginary
 210 parts (Fig. 2a), respectively. Their values are listed in Table 3. When $\lambda \geq \lambda_2$ or $\lambda \leq \lambda_1$, only purely
 211 growing or damped modes without oscillatory potentials exist, suggested by the corresponding
 212 positive or negative real parts (Fig. 2b). When $\lambda_1 < \lambda < \lambda_2$, the systems exhibit oscillatory behavior
 213 because of the presence of the imaginary parts (Fig. 2a). With the increase of λ , the models have the
 214 tendency to change from a damped oscillation to a growing oscillation. In comparison with the 4S and
 215 3S models of LY22, it appears that including the temperature equations in the system has at least two
 216 consequences:

- 217 (a) An acceleration of the oscillation, evidenced by the larger imaginary parts in the 4TS and 3TS
 218 models (Fig. 2a, orange curves) than in the 4S and 3S models (Fig. 2a, black curves).
- 219 (b) An overall stabilization for the system, evidenced by the higher λ_c in the 4TS and 3TS models
 220 than in the 4S and 3S models listed in Table 3, and the smaller real parts in the 4TS and 3TS
 221 models (Fig. 2c, orange lines) than in the 4S and 3S models (Fig. 2c, black lines).



222

223 FIG. 2. Dependences of (a) imaginary parts and (b) real parts of the oscillatory mode on λ in the 4TS (solid
 224 orange curves), 3TS (dashed orange curves), 4S (solid black curves), and 3S (dashed black curves) models. (c)
 225 is the magnified version of (b) near line $y = 0$. Results of the 4S and 3S models are from LY22. The units of
 226 the ordinate are 10^{-10} s^{-1} . The values of the other parameters are the same as those listed in Table 1. The vertical
 227 dashed gray line denotes the situation under the standard value $\lambda = 12 \text{ Sv} \cdot \text{kg}^{-1} \text{m}^3$.

228

229 TABLE 3. Values for λ_c , λ_1 and λ_2 (units: $Sv \cdot kg^{-1} m^3$) in different box models.

	4TS	3TS	4S	3S
λ_c	13.20	13.06	11.44	12.39
λ_1, λ_2	-0.92, 26.80	-0.70, 23.24	-0.89, 20.44	-0.69, 21.46

230

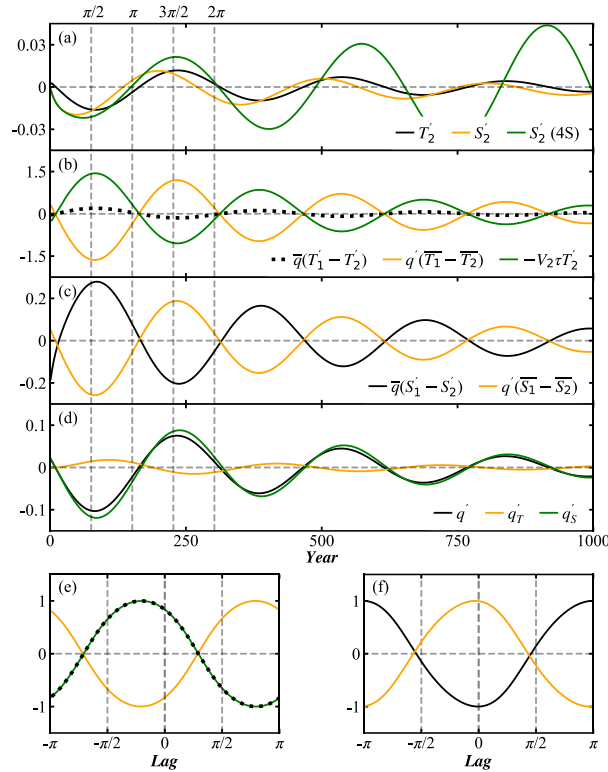
231 The stability analyses provide us the mathematical fundamentals, showing how the oscillatory
 232 behaviors of the system change when the thermal process is included. Physical insight into why the
 233 oscillation changes will be deliberated next.

234

235 3. Effects of temperature equations

236 a. Temperature feedbacks

237 There are mainly two feedbacks between the thermal process and the AMOC: the negative
 238 temperature-advection feedback and the positive restoring-advection feedback. Let us illustrate them
 239 using box 2 (Fig. 3b). Starting with a positive perturbation of q' , the anomalous advection
 240 $q'(\overline{T}_1 - \overline{T}_2)$ transports more warm water northward, T'_2 is increased and thus $\Delta\rho'_T$ is lowered, causing
 241 a decrease in q'_T . This is the negative temperature-advection feedback, which can be further illustrated
 242 by the lead/lag correlation between $q'(\overline{T}_1 - \overline{T}_2)$ and q'_T : the former leads the latter by about $\pi/4$ with
 243 a negative correlation near 1.0 (Fig. 3e, orange curve). However, the increased T'_2 also triggers a
 244 relaxation via the anomalous restoring $-V_2\tau T'_2$, whose strength is proportional to the restoring
 245 coefficient τ . This limits the growth of the positive T'_2 itself, bounding the decreases of $\Delta\rho'_T$ and q'_T .
 246 This is the positive restoring-advection feedback, which is also illustrated clearly in Fig. 3e (green
 247 curve): $-V_2\tau T'_2$ leads q'_T by about $\pi/4$ with a positive correlation near 1.0. These two feedbacks are
 248 local ones, which have the comparable amplitude and can offset each other mostly (Fig. 3b). There is
 249 a third feedback coming from $\overline{q}(T'_1 - T'_2)$ of Eq. (5b). This term is related to the remote response T'_1 ,
 250 and is a weak positive feedback (Figs. 3b, e).



251

252 FIG. 3. (a) Damped and growing oscillations in the 4TS and 4S models using the standard parameters in Table
 253 1. Black, orange and green curves are the time series of T'_2 (units: $^{\circ}\text{C}$), S'_2 (units: psu) in the 4TS model and S'_2
 254 in the 4S model, respectively. (b) Time series for temperature terms (units: $Sv \cdot ^{\circ}\text{C}$) on the right-hand side of
 255 Eq. (5b); (c) time series for salinity terms (units: $Sv \cdot psu$) on the right-hand side of Eq. (5f); (d) time series for
 256 q' , q'_T and q'_S (units: Sv) in the 4TS model. The vertical dashed gray lines in (a)-(d) mark the locations of $\pi/2$,
 257 π , $3\pi/2$, and 2π of the period (302 years) in the 4TS model. (e) Lead/lag correlation coefficients between q'_T
 258 and $\bar{q}(T'_1 - T'_2)$ (dotted black curve), $q'(\bar{T}_1 - \bar{T}_2)$ (solid orange curve) and $-V_2 \tau T'_2$ (solid green curve) in the
 259 4TS model. (f) Lead/lag correlation coefficients between q'_S and $\bar{q}(S'_1 - S'_2)$ (solid black curve), and between
 260 q'_S and $q'(\bar{S}_1 - \bar{S}_2)$ (solid orange curve) in the 4TS model. In (e)-(f), the negative lag represents q' lags the
 261 other terms.

262

263 The salt-advection feedback in the 4TS model is nearly identical to that in the 4S model of
 264 LY22. The positive and negative feedbacks come from terms $q'(\bar{S}_1 - \bar{S}_2)$ and $\bar{q}(S'_1 - S'_2)$ (Figs. 3c,
 265 f), respectively. Note that q' is the sum of salinity-induced q'_S and temperature-induced q'_T . These two
 266 components are roughly out of phase; and the former is much bigger than the latter (Fig. 3d),
 267 suggesting that the salt-advection feedback has more remarkable effect on the AMOC than the
 268 temperature-advection feedback does. In addition, although introducing the fast-restoring processes
 269 leads to an obvious time lag between the thermal process and the AMOC (Fig. 3e), there is almost no

270 time lag between the saline process and the AMOC (Fig. 3f), further suggesting the deterministic role
 271 of the saline process in the multicentennial oscillation of the AMOC.

272 Results shown in Fig. 3 are obtained from forward numerical integration of Eq. (5) with the
 273 standard parameters in Table 1. The fourth-order Runge-Kutta method is used to solve Eq. (5), with
 274 $S'_1(t = 0) = -0.02 \text{ psu}$ at the first time step, and $S' = T' = 0$. The integration time step is 7.2 days,
 275 and the total integration length exceeds 10000 years, but only the very front parts are shown.
 276 Throughout this paper, the same numerical method is employed for all experiments; and annual mean
 277 data is used for analysis.

278

279 *b. Role of restoring feedback*

280 It is the restoring feedback in the temperature equations that changes the oscillatory behavior of
 281 the system. To understand this better, let us first examine how the restoring timescale affects the
 282 temperature-advection feedback. Based on Eq. (2a), we have,

$$283 \quad \overline{T}_1 - \overline{T}_2 = (T_1^* - T_2^*) / \left[\frac{\overline{q}(V_1 + V_2)}{V_1 V_2 \tau} + 1 \right] \quad (8)$$

284 This depicts a larger τ (or a shorter restoring timescale) causes a larger $\overline{T}_1 - \overline{T}_2$, thus stronger
 285 advection $q'(\overline{T}_1 - \overline{T}_2)$. However, since the negative temperature-advection feedback is realized
 286 through an increase in T'_2 , which is in turn limited by stronger $-V_2 \tau T'_2$, the restoring-advection
 287 feedback tends to always offset mostly the temperature-advection feedback, regardless of the
 288 restoring strength (Fig. 3b).

289 TABLE 4. Conjugate eigenmodes in the 4TS model under different τ .

τ	Eigenvalues (10^{-10} s^{-1})	E-folding time (Years)	Period (Years)
$\tau_0 = 0$	$0.31 \pm 5.83i$	1025	341
$\tau_1 = (5 \text{ year})^{-1}$	$-1.48 \pm 7.30i$	-215	273
$\tau_2 = (1 \text{ year})^{-1}$	$-0.55 \pm 6.59i$	-576	302
$\tau_3 = (0.25 \text{ year})^{-1}$	$0.036 \pm 6.09i$	8830	327

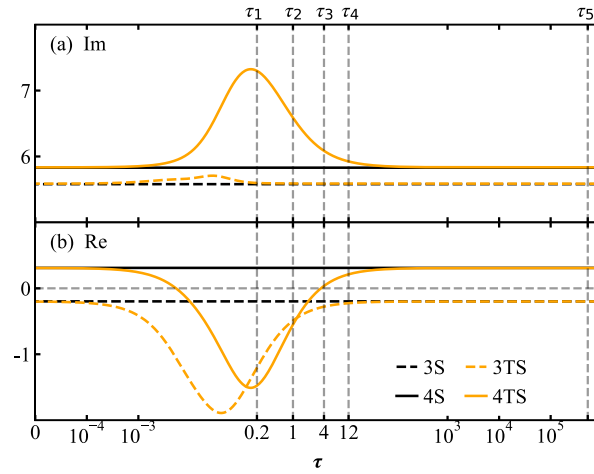
$\tau_4 = (1 \text{ month})^{-1}$	$0.21 \pm 5.92i$	1492	336
$\tau_5 = (1 \text{ minute})^{-1}$	$0.31 \pm 5.83i$	1025	341

290

291 There are two extreme situations. When $\tau \rightarrow 0$ or $\tau \rightarrow \infty$, that is, the restoring timescale for SST
 292 goes to either infinity or zero, the oscillatory eigenmode ($0.31 \pm 5.83i$) in the 4TS model is identical
 293 to that in the 4S model (Table 4). Under these two situations, the thermal process has no effect on the
 294 AMOC oscillation, and the 4TS model is practically reduced to the 4S model. In the situation with
 295 $\tau \rightarrow 0$, the linearized temperature equations (Eqs. 5a-d) are identical to the linearized salinity
 296 equations (Eqs. 5e-h), and the combined temperature and salinity equations are equivalent to the
 297 salinity equations in the 4S model. Now, $\overline{T}_1 = \overline{T}_2 = \overline{T}_3 = \overline{T}_4$ based on Eqs. (2) and (8). There is no
 298 temperature-advection feedback ($q'(\overline{T}_1 - \overline{T}_2) = 0$) anywhere, so the system is totally controlled by
 299 the saline process. In the situation with $\tau \rightarrow \infty$, $\overline{T}_1 - \overline{T}_2 = T_1^* - T_2^*$. The extremely strong restoring
 300 kills any temperature perturbations immediately, which makes $T_1' = T_2' = 0$ and the 4TS system is
 301 equivalent to a system without active thermal process, so that only the saline process matters to the
 302 oscillatory behavior. In summary, under these two extreme situations, the results from linear stability
 303 analysis suggest the oscillations of salinity and AMOC are identical to those of the 4S model in
 304 LY22; as a result, the corresponding figures (Figs. 2-3 when $\tau \rightarrow 0$ or $\tau \rightarrow \infty$) are not shown.

305 The dependences of imaginary and real parts of the oscillatory mode on τ in the 4TS model are
 306 shown in Fig. 4 (solid orange curves). Under a reasonable range of τ (from τ_1 to τ_4 ; Table 4), the
 307 oscillatory behavior of the 4TS model can vary from a damped oscillation to a weakly growing
 308 oscillation (Fig. 4b, solid orange curve). This is because the positive restoring-advection feedback
 309 becomes stronger as the restoring timescale gets shorter; and the system changes from under-
 310 compensation to overcompensation, to the negative temperature-advection feedback. Compared with
 311 the 4S model (Fig. 4b, solid black line), the 4TS model is generally more stable, manifested by the
 312 negative or longer positive e-folding time. Even for a very short SST restoring timescale (one to
 313 several months) (Table 4), the positive e-folding time of the oscillatory mode in the 4TS model is still
 314 much longer than that in the 4S model. This is because for any given τ and λ , temperature-induced q'_T
 315 is always opposite to salinity-induced q'_S , so that the total q' is always smaller, i.e., the AMOC
 316 sensitivity to buoyancy perturbation is always weaker in the 4TS model than in the 4S model. In
 317 addition, including the fast thermal restoring process leads to a shorter oscillation period in the 4TS
 318 model than in the 4S model (Fig. 4a), because the superimposition of a quick timescale and a slow

319 timescale leads to a timescale in between. Practically, since a reasonable restoring timescale is always
 320 much shorter than the multicentennial timescale, the effect of restoring timescale on the oscillation
 321 period of the system can be neglected.

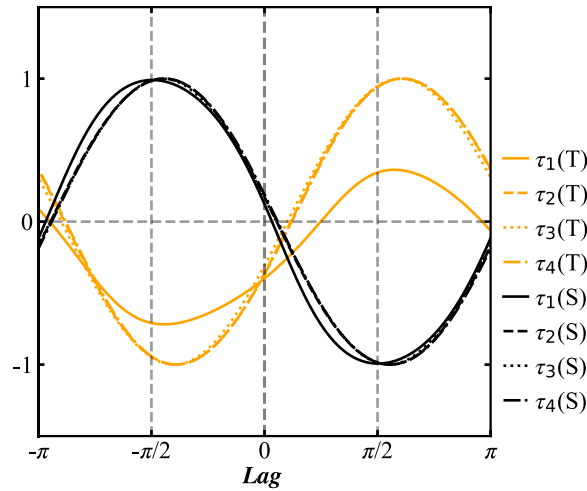


322

323 FIG. 4. Dependence of (a) positive imaginary parts and (b) real parts of the oscillatory mode on τ (units: year⁻¹)
 324 in the 4TS model (solid orange curve) and the 3TS model (dashed orange curve). The units of the ordinate are
 325 10^{-10} s^{-1} . The vertical dashed gray lines from left to right denote the situations under τ_1 , τ_2 , τ_3 , τ_4 , and τ_5 ,
 326 respectively. The reference oscillatory modes in the 4S and 3S models are plotted as the solid and dashed black
 327 lines, respectively, which are independent of τ . Here, λ is set to $12 \text{ Sv} \cdot \text{kg}^{-1} \text{ m}^3$. The values of the other
 328 parameters are the same as those listed in Table 1.

329

330 The restoring timescale also affects the relative stability of the 4TS and 3TS models. As shown
 331 in Fig. 1, under extremely strong vertical mixing in the subpolar ocean, the 4-box model (Fig. 1a) can
 332 be reduced to a 3-box model (Fig. 1b). LY22 showed that the 3S model is always more stable than the
 333 4S model. Here, we find that including the thermal process, the change of stability from the 4TS to
 334 3TS model is not that obvious (Fig. 4). To better understand the stability change, we should first
 335 recognize that whether the temperature and salinity anomalies stay in the subpolar upper or lower
 336 ocean does not influence the meridional density gradient due to the vertically weighted volume-
 337 averaged treatment. However, the time consumed in transporting temperature and salinity anomalies
 338 from the upper to lower ocean is omitted; consequently, they are removed faster from the subpolar
 339 region in the 3-box model, which reduces their restraining and amplification effects on q'_T and q'_S ,
 340 respectively. Therefore, the removals of temperature and salinity related stratifications in the 3TS
 341 model have destabilizing and stabilizing effects, respectively, on the oscillation of the system.



342

343 FIG. 5. Lead/lag correlation coefficients between $T'_2 - T'_3$ and q'_T (orange curves) and between $S'_2 - S'_3$ and q'_S
 344 (black curves) in the 4TS model under different τ . The negative lag represents q' lags the other terms. Here, λ
 345 is set to $12 Sv \cdot kg^{-1} m^3$. The values of the other parameters are the same as those listed in Table 1.

346

347 In the 4TS model, the subpolar temperature stratification $T'_2 - T'_3$ leads q'_T by about $\frac{\pi}{2}$ with a
 348 negative correlation, while the subpolar salinity stratification $S'_2 - S'_3$ leads q'_S by $\frac{\pi}{2}$ with a positive
 349 correlation (Fig. 5). Moreover, at lag 0, $T'_2 - T'_3$ ($S'_2 - S'_3$) also has negative (positive) correlation with
 350 q'_T (q'_S). These correlation relationships do not rely on the temperature restoring coefficient τ . These
 351 confirm that the existences of subpolar temperature and salinity related stratifications have stabilizing
 352 and destabilizing effects, respectively, on the system. However, whether the total subpolar buoyancy
 353 stratification plays as a stabilizing or destabilizing role depends on τ . When τ lies in the range of
 354 about several years (from τ_1 to τ_2 ; Fig. 4b), the subpolar buoyancy stratification plays as a stabilizing
 355 factor since the stabilizing effect of temperature stratification overcomes the destabilizing effect of
 356 salinity stratification. When τ is too small or too large, the temperature effect becomes weaker while
 357 the salinity effect is not influenced. Hence, the temperature stratification no longer overcomes the
 358 salinity stratification; and the 3TS model is more stable than the 4TS model, that is, including extreme
 359 mixing in the subpolar ocean can stabilize the system, as deliberated in LY22. We conclude that
 360 under realistic ranges of the parameters, the 4TS model can be more stable than the 3TS model, due to
 361 the stabilizing effect of subpolar temperature stratification.

362

363 4. Realization of self-sustained oscillation

364 Self-sustained oscillation is still absent in the 4TS model. Under the same parameters, the 4TS
 365 model is more stable than the 4S model in LY22 (Fig. 2c), as discussed in section 3. However, this
 366 does not lead to a self-sustained oscillation in the 4TS model; clearly, additional processes are
 367 needed. In LY22, an enhanced mixing process is added explicitly in the subpolar ocean, to realize a
 368 self-sustained oscillation. There is also an alternative way to realize a self-sustained oscillation as
 369 shown in Rivin and Tziperman (1997) (hereafter RT97), in which a nonlinear relationship between
 370 the AMOC strength and meridional density gradient is employed. Here, we want to emphasize that a
 371 self-sustained oscillation should first satisfy the instability criterion detailed in LY22, that is, $\lambda > \lambda_c$,
 372 depicting that the AMOC should be sensitive enough to the perturbation of meridional density
 373 gradient, and the intrinsic oscillatory mode is an unstable mode. When $\lambda \leq \lambda_c$, the oscillatory mode
 374 itself is a decayed or neutral mode; and any additional mixing or nonlinear processes will make the
 375 oscillation even more decayed.

376

377 *a. Self-sustained oscillation with enhanced subpolar mixing*

378 Similar to LY22, we introduce an enhanced mixing term between boxes 2 and 3 in the 4TS
 379 model. Eqs. (5b-c) and (5f-g) become,

$$380 \quad V_2 \dot{T}'_2 = q'(\overline{T}_1 - \overline{T}_2) + \overline{q}(T'_1 - T'_2) - k_m(T'_2 - T'_3) - V_2 \tau T'_2 \quad (9a)$$

$$381 \quad V_3 \dot{T}'_3 = \overline{q}(T'_2 - T'_3) + k_m(T'_2 - T'_3) \quad (9b)$$

$$382 \quad V_2 \dot{S}'_2 = q'(\overline{S}_1 - \overline{S}_2) + \overline{q}(S'_1 - S'_2) - k_m(S'_2 - S'_3) \quad (9c)$$

$$383 \quad V_3 \dot{S}'_3 = \overline{q}(S'_2 - S'_3) + k_m(S'_2 - S'_3) \quad (9d)$$

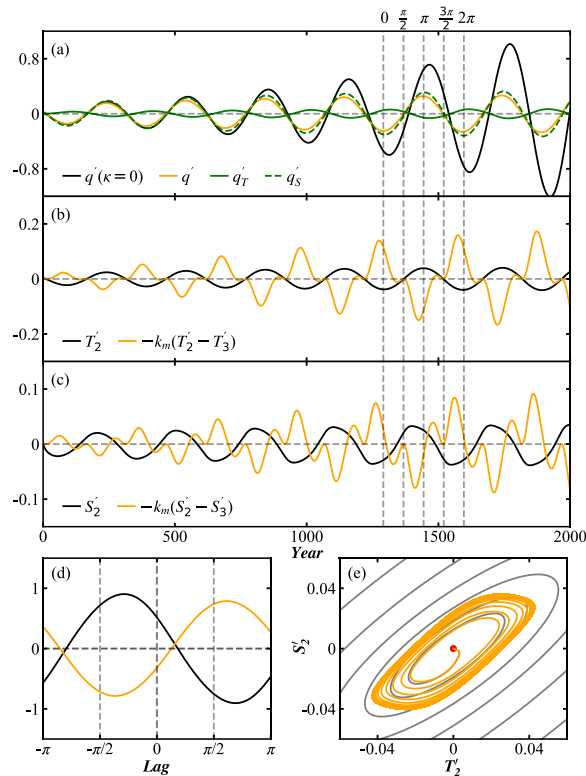
384 And the mixing coefficient k_m (units: m^3/s) is represented by:

$$385 \quad k_m = \kappa q'^2 \quad (9e)$$

386 Here, κ (units: $m^{-3}s$) is a positive constant. We set it to $1 \times 10^{-4} m^{-3}s$ in this paper. No matter the
 387 sign of q' , k_m is always positive and helps remove the subpolar upper-ocean anomalies. Detailed
 388 physics of the enhanced mixing process was discussed in LY22.

389 A growing oscillation (Fig. 6a, solid black curve) is turned into a self-sustained oscillation (Fig.
 390 6a, solid orange curve) when enhanced subpolar mixing is included. Here, $\lambda = 14 Sv \cdot kg^{-1} m^3$; and
 391 the intrinsic mode of the 4TS model is unstable. As q' grows (decreases), more warm and saline (cold

392 and fresh) water is removed from the subpolar upper ocean, which enters the lower ocean through
 393 anomalous mixings $-k_m(T'_2 - T'_3)$ and $-k_m(S'_2 - S'_3)$ (Figs. 6b, c, solid orange curves). In turn,
 394 further growth (decrease) of q' is restrained. Beware that the temperature and salinity mixing-
 395 advection feedbacks have destabilizing and stabilizing effects on q' , respectively (Fig. 6d). Their
 396 combined effect on subpolar density is to stabilize q' . In summary, including enhanced mixing in the
 397 subpolar ocean can well establish a self-sustained oscillation, which can be seen more clearly in the
 398 phase diagram of T'_2 vs S'_2 (Fig. 6e, orange curve); that is, a limit cycle is formed eventually.



399

400 FIG. 6. Oscillations under $\lambda = 14 \text{ Sv} \cdot \text{kg}^{-1} \text{ m}^3$. (a) Time series for q' (solid black curve) under $\kappa = 0$, q'
 401 (solid orange curve), q_T' (solid green curve) and q_S' (dashed green curve) under $\kappa = 1 \times 10^{-4} \text{ m}^{-3} \text{ s}$ (units:
 402 Sv). (b) Time series for T'_2 (solid black curve; units: $^{\circ}\text{C}$) and $-k_m(T'_2 - T'_3)$ (solid orange curve; units: $\text{Sv} \cdot ^{\circ}\text{C}$).
 403 (c) Time series for S'_2 (solid black curve; units: psu) and $-k_m(S'_2 - S'_3)$ (solid orange curve; units: $\text{Sv} \cdot \text{psu}$).
 404 (d) Lead/lag correlation coefficients for $-k_m(T'_2 - T'_3)$ and q_T' (solid black curve), $-k_m(S'_2 - S'_3)$ and q_S'
 405 (solid orange curve). (e) T'_2 - S'_2 phase space diagrams for years 1-10000. The red dot represents the initial
 406 location of T'_2 and S'_2 . Black curve is for $\kappa = 0$, and orange curve, for $\kappa = 1 \times 10^{-4} \text{ m}^{-3} \text{ s}$. The vertical dashed
 407 gray line in (a), (b) and (c) marks a holonomic oscillation period under $\kappa = 1 \times 10^{-4} \text{ m}^{-3} \text{ s}$. The values of the
 408 other parameters are the same as those listed in Table 1.
 409

410 The oscillation period under $\kappa = 1 \times 10^{-4} \text{ m}^{-3} \text{ s}$ is 300 years, quite close to the 306-year
 411 analytical period under $\kappa = 0$. The period is insensitive to the value of κ , indicating that it is

412 dominated by advection instead of the mixing process. The κ chosen here is one order smaller than
 413 the value used in LY22, suggesting that even a small bounding from subpolar vertical mixing can lead
 414 to self-sustained oscillation.

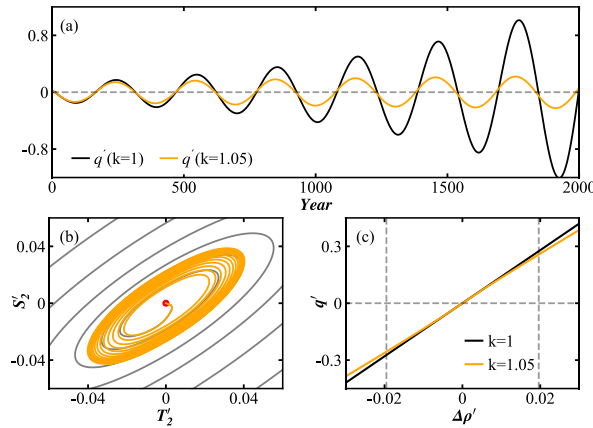
415

416 *b. Self-sustained oscillation with nonlinear AMOC-density relation*

417 It has been prevailing to set the sensitivity of the AMOC in low-order models to be linearly
 418 proportional to the meridional density gradient (Stommel 1961; GT95; Zhao et al. 2016; Shi and
 419 Yang 2021). As an alternative, we use a nonlinear relation analogous to Cessi (1994) and RT97 to
 420 introduce a degree of nonlinearity in this study. Now, Eq. (4b) becomes,

$$421 \quad q' = \begin{cases} \lambda \rho_{cri} \left[k \left(\left[\frac{\Delta \rho'}{\rho_{cri}} \right]^{\frac{1}{k}} - 1 \right) + 1 \right], & \text{if } \Delta \rho' > \rho_{cri} \\ \lambda \Delta \rho' & \text{if } -\rho_{cri} < \Delta \rho' < \rho_{cri} \\ -\lambda \rho_{cri} \left[k \left(\left[-\frac{\Delta \rho'}{\rho_{cri}} \right]^{\frac{1}{k}} - 1 \right) + 1 \right], & \text{if } \Delta \rho' < -\rho_{cri} \end{cases} \quad (10)$$

422 At $k = 1$, Eq. (10) is reduced to the linear Eq. (4b); and the system exhibits growing oscillation under
 423 $\lambda = 14 Sv \cdot kg^{-1} m^3$ (Figs. 6a, 7a, black curves). If $k = 1.05$ with $\rho_{cri} = 0.002 kg/m^3$, a small
 424 degree of nonlinearity (Fig. 7c, orange curve) will be introduced into the linear system. The self-
 425 sustained oscillation is then realized (Fig. 7a, orange curve), and a limit-cycle is achieved (Fig. 7b,
 426 orange curve). The intersections between the vertical dashed gray lines and the abscissa axis in Fig.
 427 7c mark the upper and lower limits for $\Delta \rho'$ during the integration. As $\Delta \rho'$ grows, the nonlinear
 428 bounding effect of Eq. (10) gradually emerges, limiting the fluctuation tendency of q' . The bounding
 429 manifested as the difference between the solid and orange curves is very small (Fig. 7c). Hence, even
 430 a tiny degree of internal nonlinearity from the AMOC-meridional density gradient relation can lead to
 431 self-sustained oscillation. The period here is 303 years, hardly deviated from the 306-year eigen
 432 period of the linear system, reflecting again the robustness of the advection-dominated eigenmode.



433

434 FIG. 7. Oscillations under $\lambda = 14 \text{ Sv} \cdot \text{kg}^{-1} \text{ m}^3$. (a) Time series for q' (units: Sv) under $k = 1$ (black curve)
 435 and $k = 1.05$ (orange curve). (b) T_2' - S_2' phase space diagrams for years 1-10000. The red dot represents the
 436 initial location of T_2' and S_2' . Black curve is for $k = 1$, and orange curve, for $k = 1.05$. (c) Variations of q' with
 437 $\Delta\rho'$ (units: kg/m^3) under $k = 1$ (black curve) and $k = 1.05$ (orange curve). The intersections between the
 438 vertical dashed gray lines and the abscissa axis mark the upper and lower limits for $\Delta\rho'$ during the integration.
 439 The values of the other parameters are the same as those listed in Table 1.

440

441 5. Eigenmode sensitivity

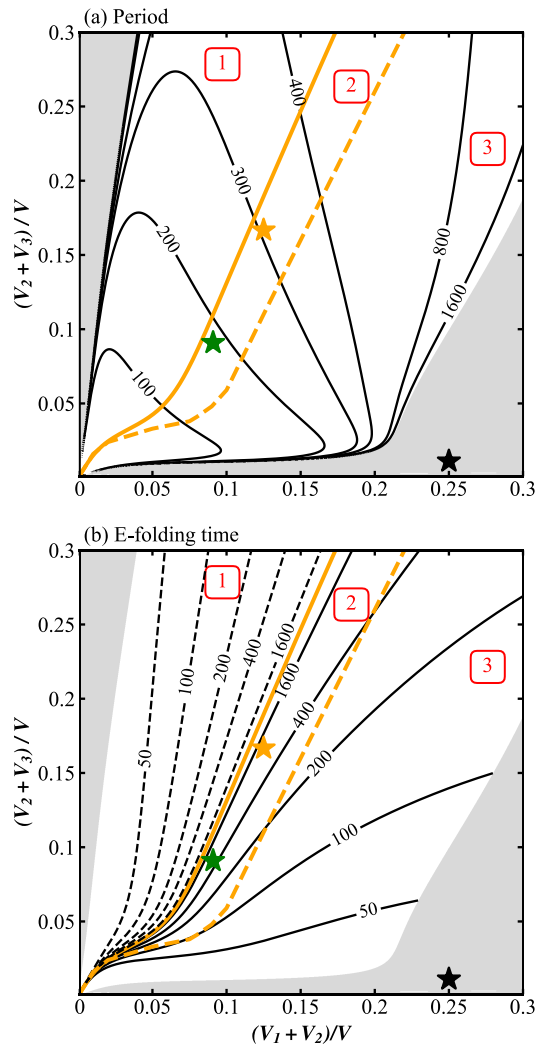
442 In theoretical models, model parameters can be tuned to control oscillation properties. GT95
 443 identified a multidecadal oscillation in their 4TS model, whose period is far shorter than our
 444 multicentennial period. RT97 used a 3TS model, and also identified a multi-decadal mode. The
 445 system stability in RT97's model is far lower than ours. Since both GT95 and RT97 also employed
 446 mixed boundary conditions, and their core dynamics are all advection feedbacks, such differences in
 447 eigenmodes are likely to originate from parameter choices. Previously studies usually tuned
 448 parameters to study multi-equilibria problems (Colin de Verdière et al. 2006; Colin de Verdière 2007;
 449 Sévellec et al. 2010; Sévellec and Fedorov 2014). In single-equilibrium oscillation studies, the
 450 parameters effects on the eigenmode have not been widely heeded, which will be addressed next via
 451 numerical stability analyses.

452

453 a. Effect of basin geometry

454 Basin geometry can affect both the period and e-folding time of the system eigenmode (Fig. 8).
 455 The eigen period increases roughly monotonously with the increases of both the subpolar ocean
 456 fraction $(V_2 + V_3)/V$ and the upper ocean fraction $(V_1 + V_2)/V$ (Fig. 8a). The standard geometry is

457 $(V_1 + V_2)/V = 1/8$ and $(V_2 + V_3)/V = 1/6$ in this paper, denoted by the orange stars in Fig. 8. In
 458 GT95, the fractions of the upper and subpolar ocean boxes are both $1/11$, falling in the lower left
 459 corner of Fig. 8 (denoted by the green star), with a period less than 200 years if \bar{q} is set to $10 Sv$.
 460 Actually, the \bar{q} in GT95 was set to a higher value of about $17 Sv$, representing a much faster
 461 overturning rate; thus, the period is further shortened to the century scale. Consequently, it is
 462 reasonable to deduce that the multi-decadal period in GT95 is not at odds with our multicentennial
 463 period. Popularity of multi-decadal phenomena back then might account for their choice of model
 464 parameters.



465

466 FIG. 8. Sensitivity of (a) period (units: years) and (b) e-folding time (units: years) of the eigenmode to subpolar
 467 ocean fraction $(V_2 + V_3)/V$ and upper ocean fraction $(V_1 + V_2)/V$ under $\lambda = 14 Sv \cdot kg^{-1} m^3$. The orange
 468 star denotes the mode with standard values used in this work. The green and black stars denote the standard
 469 values used in GT95 and RT97, respectively. The solid orange curve is both the stability threshold and the
 470 lower limit of probability for self-sustained oscillation. The dashed orange curve is the upper limit of

471 probability for self-sustained oscillation. The light gray areas correspond to purely damped or growing regime
 472 without the imaginary part. The oscillatory mode is damped in region 1, potentially self-sustained in region 2 if
 473 bounding terms are affiliated, and growing in region 3. The values of the other parameters are the same as
 474 those listed in Table 1.
 475

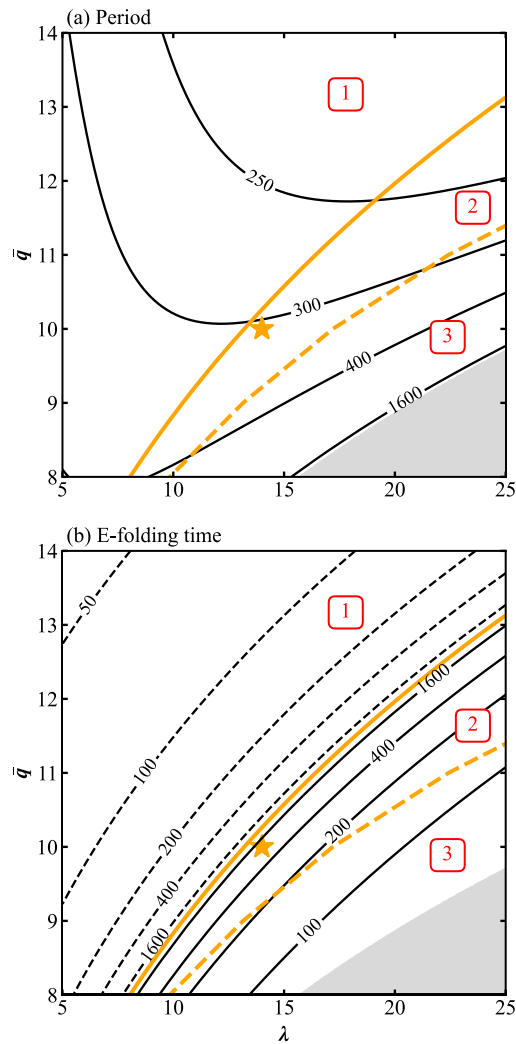
476 In Fig. 8b, a higher $(V_1 + V_2)/V$ and a lower $(V_2 + V_3)/V$ are linked to lower stability. This
 477 explains why the mode in RT97 can be easily unstable even under a low AMOC sensitivity
 478 (equivalent to $\lambda = 5.7 Sv \cdot kg^{-1} m^3$ in this paper). The basin geometry of RT97 is denoted by the
 479 black star in Fig. 8. Their high-latitude box stands for a small deep-water formation region instead of
 480 the subpolar region, so it was set to only around 1/100 the volume of the entire ocean basin. However,
 481 their upper ocean is as large as 1/4 of the entire ocean basin. Therefore, the low stability seen in RT97
 482 owes to their volume configuration, according to our stability analyses. The light gray areas in Fig. 8
 483 denote purely damped or growing region without oscillatory potentials. The solid orange curve
 484 partitions the stable and unstable regions, making itself also the lower limit for a possible self-
 485 sustained oscillation. The oscillatory mode is damped in region 1 because of negative real part, and is
 486 growing due to positive real part in regions 2 and 3 (Fig. 8b). Nevertheless, we find that only in
 487 region 2 could the self-sustained oscillation take place if bounding terms are affiliated, through a large
 488 number of numerical experiments (not shown). No self-sustained oscillation is able to exist in region
 489 3, which is separated from region 2 by the dashed orange curve. Compared to Fig. 9 in LY22, the
 490 probability for a self-sustained oscillation is increased in the 4TS model, since the area of region 2 is
 491 larger than that in the 4S model of LY22.

492

493 *b. Effect of mean flow*

494 Given the meridional density gradient, the total AMOC strength q is determined by its
 495 sensitivity λ to the meridional density gradient and the equilibrium strength \bar{q} . The period decreases
 496 monotonically as \bar{q} increases (Fig. 9a), reflecting that a faster overturning leads to a shorter oscillation
 497 period. Larger λ and smaller \bar{q} are both destabilizing factors (Fig. 9b). A larger λ results in more
 498 intense fluctuation of q' under the same perturbation of meridional density gradient, contributing to a
 499 less stable system. A decreased \bar{q} weakens the equilibrium advection terms $\bar{q}(T'_1 - T'_2)$ and
 500 $\bar{q}(S'_1 - S'_2)$, therefore limiting their destabilizing and stabilizing effects, respectively. The latter is
 501 more evident due to the dominant role that salinity plays in establishing AMOC variability. The

502 combined effect of temperature and salinity advection under a smaller \bar{q} is to make the system more
 503 unstable.



504

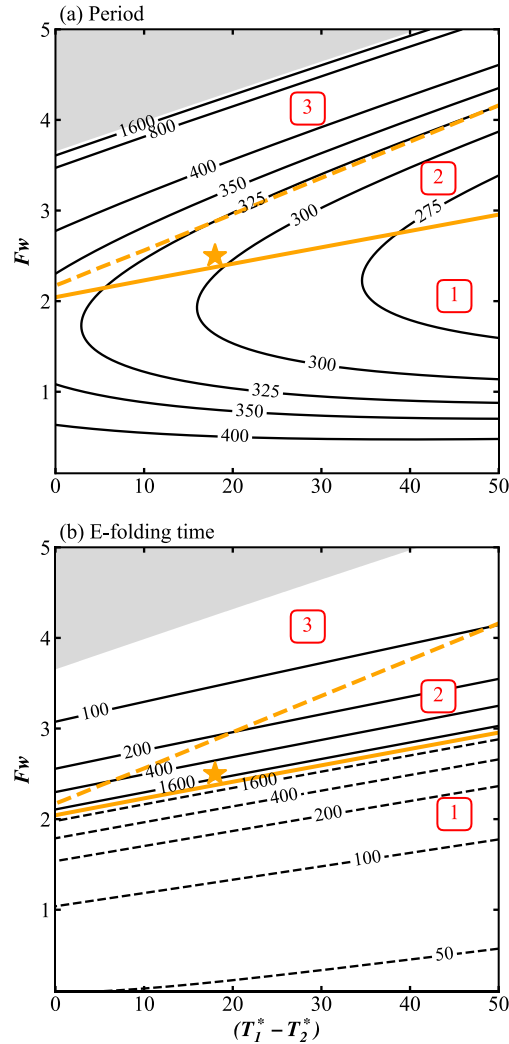
505 FIG. 9. Same as Fig. 8, but the ordinate and abscissa correspond to the equilibrium AMOC strength \bar{q} (units:
 506 Sv) and the linear closure coefficient λ (units: $Sv \cdot kg^{-1} m^3$), respectively.

507

508 *c. Effect of boundary conditions*

509 Mixed boundary conditions are adopted in the 4TS model, where T_1^* and T_2^* control the surface
 510 heat flux while F_w controls the surface virtual salt flux. The meridional restoring temperature gradient
 511 $T_1^* - T_2^*$ influences the system, while the exact values of T_1^* and T_2^* have no effect. Figure 10a shows
 512 that the period shortens marginally with the increase of $T_1^* - T_2^*$, but exhibits a decrease-to-increase
 513 tendency as F_w grows. Smaller $T_1^* - T_2^*$ and larger F_w all lead to lower stability (Fig. 10b). From Eq.
 514 (8), we derive that $\bar{T}_1 - \bar{T}_2$ lowers as $T_1^* - T_2^*$ decreases; therefore, the temperature effects are

515 hampered due to the weaker negative temperature-advection feedback. We illustrated in section 3.1
 516 that a faster restoring denoted by a smaller $1/\tau$ also limits the temperature effects. Hence, the
 517 temperature effects are promoted by the increases in $T_1^* - T_2^*$ and $1/\tau$, followed by a more stable
 518 system with a shorter period. It can also be seen from Eq. (2b) that a larger F_w increases $q'(\overline{S}_1 - \overline{S}_2)$,
 519 so the destabilizing positive salt-advection feedback is reinforced, which is consistent with the finding
 520 of Sévellec et al. (2006).



521

522 FIG. 10. Same as Fig. 8, but the ordinate and abscissa correspond to the freshwater flux F_w (units: $10 \text{ psu} \cdot \text{Sv}$)
 523 and the meridional restoring temperature gradient $T_1^* - T_2^*$ (units: $^{\circ}\text{C}$), respectively.

524

525 **6. Summary and discussion**

526 As the second part of our theoretical studies on AMOC multicentennial variability, this study
 527 complements LY22 by including temperature equations in the box model. Mixed boundary conditions

528 are employed for surface temperature and salinity. The thermal process includes the negative
529 temperature-advection feedback and positive restoring-advection feedback. The latter never overruns
530 the former; thus, the resultant temperature feedback is negative. Including the thermal process leads to
531 an acceleration of oscillation because of the fast thermal-restoring process, a stabilization force for the
532 system because of the negative temperature-advection feedback and a portion of stabilization for the
533 subpolar stratification due to temperature stratification.

534 Similar to LY22, bounding processes are needed to realize a self-sustained oscillation, which can
535 be either enhanced subpolar mixing or weak nonlinearity in the AMOC-meridional density gradient
536 relation. Multicentennial eigenmode is robust regardless of these bounding processes, because the
537 oscillatory eigenmode is fundamentally determined by advection processes. Only a tiny magnitude of
538 the bounding process is able to realize the self-sustained oscillation. Same as in LY22, the effects of
539 nonlinear temperature and salinity advection terms on the self-sustained oscillation are trivial and can
540 be safely neglected (figure not shown), and the external stochastic forcing can also excite a
541 sustainable multicentennial oscillation (figure not shown). Compared to the 4S model in LY22, the
542 probability for a self-sustained oscillation in the 4TS model is much increased with temperature
543 equations added.

544 Stability analyses reveal that the period and stability of the oscillatory eigenmode are sensitive to
545 model geometry, flow properties and boundary conditions. Generally, smaller subpolar and upper
546 oceans tend to shorten the period. Larger subpolar ocean and smaller upper ocean have stabilizing
547 effects on the system. A stronger AMOC shortens the period due to the faster overturning rate,
548 stabilizing the system through balancing the positive salt-advection feedback more quickly. Higher
549 AMOC sensitivity to the meridional density gradient makes the system less stable. Increasing surface
550 freshwater flux energizes the destabilizing salt-advection feedback, and lowers the system stability,
551 because the background meridional salinity gradient will be stronger. It also lengthens the period of
552 the system because more time is needed to consume the stronger background salinity gradient. Larger
553 meridional restoring temperature gradient strengthens the thermal process; thus, it shortens the period
554 and increases the system stability.

555 The box model is highly idealized, aimed at providing heuristic understanding of the
556 multicentennial AMOC oscillation. This work can also help us understand the prevalence of
557 centennial to multicentennial AMOC oscillations found in a few pre-industrial control runs using
558 high-order models (Vellinga and Wu 2004; Park and Latif 2008; Delworth and Zeng 2012; Yang et

559 al. 2015; Jiang et al. 2021). Whether the multicentennial AMOC oscillations found in Earth system
560 models are self-sustained or stochastically-sustained is obscured by their intricate model physics.
561 However, our study suggests that a self-sustained oscillation can appear as long as a tiny magnitude
562 of nonlinearity or additional mixing is included, which is easy to realize in the realistic ocean.
563 Thereby, we conclude that even with random components removed completely, self-sustained
564 multicentennial oscillation has a good chance to exist in high-order models.

565 The core of the oscillation mechanism here is the advection process, consistent with many
566 previous studies (Mikolajewicz and Maier-Reimer 1990; Winton and Sarachik 1993; Drijfhout et al.
567 1996; Delworth and Zeng 2012). Sensitivity of period to model geometry is observed not only in our
568 theoretical model but also in higher complexity models (Weaver and Sarachik 1991; Drijfhout et al.
569 1996; Delworth and Zeng 2012). Since the flow rate and route affect the overturning rate, a more
570 precise simulation of AMOC structure and a finer model resolution are likely to improve the
571 simulation of AMOC oscillation. Although the boundary conditions in our model influence the
572 eigenmode, the essence for such impact is climate feedbacks. It inspires us that a better representation
573 of climate feedbacks in high-order models may improve their performances.

574 The warming and freshwater hosing in the North Atlantic will reduce the meridional temperature
575 gradient while enhance the meridional salinity gradient, hampering the negative temperature-
576 advection feedback and strengthening the positive salt-advection feedback. On one hand, this implies
577 that the AMOC might march gradually toward (not necessarily reach) the collapse state (Gregory et
578 al. 2005; Sévellec et al. 2017; Dai 2022), since its stability is likely to reduce, as revealed in this
579 paper. On the other hand, this also implies that the period for the multicentennial timescale portion of
580 the AMOC oscillation is likely to be lengthened in the future, which has not gained attention yet.
581 However, the portion with decadal to multi-decadal periods of the AMOC is believed to be shortened
582 under global warming scenario based on Rossby wave dynamics (Cheng et al. 2016; Ma et al. 2021).
583 As global warming persists, more attention should be paid to how the multicentennial AMOC period
584 would change in the future, since the global warming might occur on the background of a
585 multicentennial oscillation.

586 This theoretical study can be improved in several aspects. The one-hemisphere configuration
587 singles out only North Atlantic advection, and contributions from other ocean basins are not
588 considered. Extending the one-hemisphere model into an inter-hemisphere one as in Scott et al.
589 (1999) and in Lucarini and Stone (2005), or incorporating the Arctic Ocean as in Lambert et al.

590 (2016) may provide more insightful results. Too few natural feedbacks are reserved in our model
591 because of the ocean-only configuration and mixed boundary conditions. Adding more feedbacks,
592 such as meridional moisture transport feedback (Tziperman and Gildor 2002), wind forcing feedback
593 (Sherriff-Tadano and Abe-Ouchi 2020) and sea ice feedback (Jayne and Marotzke 1999), should
594 improve the authenticity of stability and other characteristics of AMOC oscillation.

595

596 *Acknowledgement:* This research is jointly supported by the NSF of China (Nos. 42230403, 41725021
597 and 91737204) and by the foundation at the Shanghai Frontiers Science Centre of Atmosphere-Ocean
598 Interaction of Fudan University.

599

600 **Data Availability Statement:**

601 This is a theory-based article and no datasets were generated during the current study.

602

References

- 603
- 604 Bretherton, F. P., 1982: Ocean climate modeling. *Prog. Oceanogr.*, **11**, 93-129.
- 605 Cessi, P., 1994: A simple box model of stochastically forced thermohaline flow. *J. Phys. Oceanogr.*,
- 606 **24**, 1911-1920.
- 607 Chabaud, L., M. F. S. Goni, S. Desprat, and L. Rossignol, 2014: Land-sea climatic variability in the
- 608 eastern north Atlantic subtropical region over the last 14,200 years: atmospheric and oceanic
- 609 processes at different timescales. *Holocene*, **24**, 787-797.
- 610 Cheng, J., Z. Y. Liu, S. Q. Zhang, W. Liu, L. N. Dong, P. Liu, and H. L. Li, 2016: Reduced
- 611 interdecadal variability of Atlantic meridional overturning circulation under global warming.
- 612 *Proc. Natl. Acad. Sci. USA*, **113**, 3175-3178.
- 613 Cimadoribus, A. A., S. S. Drijfhout, and H. A. Dijkstra, 2014: Meridional overturning circulation:
- 614 stability and ocean feedbacks in a box model. *Climate Dyn.*, **42**, 311-328.
- 615 Colin de Verdière, A., 2007: A simple model of millennial oscillations of the thermohaline circulation.
- 616 *J. Phys. Oceanogr.*, **37**, 1142-1155.
- 617 Colin de Verdière, A., 2010: The instability of the thermohaline circulation in a low-order model. *J.*
- 618 *Phys. Oceanogr.*, **40**, 757-773.
- 619 Colin de Verdière, A., M. Ben Jelloul, and F. Sévellec, 2006: Bifurcation structure of thermohaline
- 620 millennial oscillations. *J. Climate*, **19**, 5777-5795.
- 621 Dai, A. G., 2022: Arctic amplification is the main cause of the Atlantic meridional overturning
- 622 circulation weakening under large CO₂ increases. *Climate Dyn.*, **58**, 3243-3259.
- 623 Delworth, T. L., and F. R. Zeng, 2012: Multicentennial variability of the Atlantic meridional
- 624 overturning circulation and its climatic influence in a 4000 year simulation of the GFDL CM2.1
- 625 climate model. *Geophys. Res. Lett.*, **39**.
- 626 Drijfhout, S., C. Heinze, M. Latif, and E. MaierReimer, 1996: Mean circulation and internal
- 627 variability in an ocean primitive equation model. *J. Phys. Oceanogr.*, **26**, 559-580.
- 628 Gregory, J. M., and Coauthors, 2005: A model intercomparison of changes in the Atlantic
- 629 thermohaline circulation in response to increasing atmospheric CO₂ concentration. *Geophys.*
- 630 *Res. Lett.*, **32**, 5.
- 631 Griffies, S. M., and E. Tziperman, 1995: A linear thermohaline oscillator driven by stochastic
- 632 atmospheric forcing. *J. Climate*, **8**, 2440-2453.
- 633 Haney, R. L., 1971: Surface thermal boundary condition for ocean circulation models. *J. Phys.*
- 634 *Oceanogr.*, **1**, 241-248.
- 635 Huang, R. X., J. R. Luyten, and H. M. Stommel, 1992: Multiple equilibrium states in combined

- 636 thermal and saline circulation. *J. Phys. Oceanogr.*, **22**, 231-246.
- 637 Jayne, S. R., and J. Marotzke, 1999: A destabilizing thermohaline circulation-atmosphere-sea ice
638 feedback. *J. Climate*, **12**, 642-651.
- 639 Jiang, W. M., G. Gastineau, and F. Codron, 2021: Multicentennial variability driven by salinity
640 exchanges between the Atlantic and the Arctic ocean in a coupled climate model. *J. Adv. Model.*
641 *Earth Syst.*, **13**, e2020MS002366.
- 642 Joyce, T. M., 1991: Thermohaline catastrophe in a simple four-box model of the ocean climate.
643 *Journal of Geophysical Research-Oceans*, **96**, 20393-20402.
- 644 Lambert, E., T. Eldevik, and P. M. Haugan, 2016: How northern freshwater input can stabilise
645 thermohaline circulation. *Tellus Ser. A-Dyn. Meteorol. Oceanol.*, **68**, 15.
- 646 Li, Y., and H. Yang, 2022: A theory for self-sustained multicentennial oscillation of the Atlantic
647 meridional overturning circulation. *J. Climate*, **35**, 5883-5896.
- 648 Lucarini, V., and P. H. Stone, 2005: Thermohaline circulation stability: A box model study. Part I:
649 Uncoupled model. *J. Climate*, **18**, 501-513.
- 650 Ma, X. F., W. Liu, N. J. Burls, C. L. Chen, J. Cheng, G. Huang, and X. C. Li, 2021: Evolving AMOC
651 multidecadal variability under different CO2 forcings. *Climate Dyn.*, **57**, 593-610.
- 652 Marotzke, J., 1996: Analysis of thermohaline feedbacks, in Decadal Climate Variability: Dynamics
653 and Predictability, D. L. T. Anderson and J. Willebrand, eds., NATO ASI Series, Series I, 333-
654 378.
- 655 Marotzke, J., and P. H. Stone, 1995: Atmospheric transports, the thermohaline circulation, and flux
656 adjustments in a simple coupled model. *J. Phys. Oceanogr.*, **25**, 1350-1364.
- 657 Mikolajewicz, U., and E. Maier-Reimer, 1990: Internal secular variability in an ocean general
658 circulation model. *Climate Dyn.*, **4**, 145-156.
- 659 Muir, L. C., and A. V. Fedorov, 2015: How the AMOC affects ocean temperatures on decadal to
660 centennial timescales: the North Atlantic versus an interhemispheric seesaw. *Climate Dyn.*, **45**,
661 151-160.
- 662 Park, W., and M. Latif, 2008: Multidecadal and multicentennial variability of the meridional
663 overturning circulation. *Geophys. Res. Lett.*, **35**.
- 664 Pierce, D. W., 1996: Reducing phase and amplitude errors in restoring boundary conditions. *J. Phys.*
665 *Oceanogr.*, **26**, 1552-1560.
- 666 Rahmstorf, S., 1996: On the freshwater forcing and transport of the Atlantic thermohaline circulation.
667 *Climate Dyn.*, **12**, 799-811.
- 668 Rahmstorf, S., and J. Willebrand, 1995: The role of temperature feedback in stabilizing the

- 669 thermohaline circulation. *J. Phys. Oceanogr.*, **25**, 787-805.
- 670 Rivin, I., and E. Tziperman, 1997: Linear versus self-sustained interdecadal thermohaline variability
671 in a coupled box model. *J. Phys. Oceanogr.*, **27**, 1216-1232.
- 672 Roebber, P. J., 1995: Climate variability in a low-order coupled atmosphere-ocean model. *Tellus Ser.*
673 *A-Dyn. Meteorol. Oceanol.*, **47**, 473-494.
- 674 Schmidt, G. A., and L. A. Mysak, 1996: The stability of a zonally averaged thermohaline circulation
675 model. *Tellus Ser. A-Dyn. Meteorol. Oceanol.*, **48**, 158-178.
- 676 Scott, J. R., J. Marotzke, and P. H. Stone, 1999: Interhemispheric thermohaline circulation in a
677 coupled box model. *J. Phys. Oceanogr.*, **29**, 351-365.
- 678 Sévellec, F., and A. V. Fedorov, 2014: Millennial variability in an idealized ocean model: predicting
679 the AMOC regime shifts. *J. Climate*, **27**, 3551-3564.
- 680 Sévellec, F., T. Huck, and M. Ben Jelloul, 2006: On the mechanism of centennial thermohaline
681 oscillations. *J. Mar. Res.*, **64**, 355-392.
- 682 Sévellec, F., T. Huck, and C. d. V. A., 2010: From centennial to millennial oscillation of the
683 thermohaline circulation. *J. Mar. Res.*, **68**, 723-742.
- 684 Sévellec, F., A. V. Fedorov, and W. Liu, 2017: Arctic sea-ice decline weakens the Atlantic meridional
685 overturning circulation. *Nat. Climate Change*, **7**, 604-610.
- 686 Sherriff-Tadano, S., and A. Abe-Ouchi, 2020: Roles of sea ice-surface wind feedback in maintaining
687 the glacial Atlantic meridional overturning circulation and climate. *J. Climate*, **33**, 3001-3018.
- 688 Shi, J. Q., and H. J. Yang, 2021: Bjerknes compensation in a coupled global box model. *Climate Dyn.*,
689 **57**, 3569-3582.
- 690 Stommel, H., 1961: Thermohaline convection with two stable regimes of flow. *Tellus*, **13**, 224-230.
- 691 Tziperman, E., and H. Gildor, 2002: The stabilization of the thermohaline circulation by the
692 temperature-precipitation feedback. *J. Phys. Oceanogr.*, **32**, 2707-2714.
- 693 Vellinga, M., and P. L. Wu, 2004: Low-latitude freshwater influence on centennial variability of the
694 Atlantic thermohaline circulation. *J. Climate*, **17**, 4498-4511.
- 695 Walin, G., 1985: The thermohaline circulation and the control of ice ages. *Palaeogeogr. Palaeocl.*
696 *Palaeoecol.*, **50**, 323-332.
- 697 Weaver, A. J., and E. S. Sarachik, 1991: Evidence for decadal variability in an ocean general-
698 circulation model - an advective mechanism. *Atmos. Ocean*, **29**, 197-231.
- 699 Winton, M., and E. S. Sarachik, 1993: Thermohaline oscillations induced by strong steady salinity
700 forcing of ocean general-circulation models. *J. Phys. Oceanogr.*, **23**, 1389-1410.
- 701 Yang, H. J., Q. Li, K. Wang, Y. Sun, and D. X. Sun, 2015: Decomposing the meridional heat transport

- 702 in the climate system. *Climate Dyn.*, **44**, 2751-2768.
- 703 Zhang, R., and Coauthors, 2019: A review of the role of the Atlantic meridional overturning
704 circulation in Atlantic multidecadal variability and associated climate impacts. *Rev. Geophys.*,
705 **57**, 316-375.
- 706 Zhang, S., R. J. Greatbatch, and C. A. Lin, 1993: A reexamination of the polar halocline catastrophe
707 and implications for coupled ocean atmosphere modeling. *J. Phys. Oceanogr.*, **23**, 287-299.
- 708 Zhao, Y. Y., H. J. Yang, and Z. Y. Liu, 2016: Assessing Bjerknes compensation for climate variability
709 and its time-scale dependence. *J. Climate*, **29**, 5501-5512.
- 710