



1 An Earth system deep learning classifier for tipping point 2 detection

3 Madleen Grohganz¹, Thomas M. Bury², Bregje van der Bolt³, Gert-Jan Reichart^{1,4}, Rick
4 Hennekam¹

5 ¹Department of Ocean Systems, NIOZ Royal Netherlands Institute for Sea Research, Den Burg, 1790 AB, The
6 Netherlands

7 ² Department of Mathematics, University of California, Riverside, Riverside, CA 92521, United States

8 ³Environmental Sciences Group, Water Systems and Global Change, Wageningen University, Wageningen, 6708 PB,
9 The Netherlands

10 ⁴Department of Earth Sciences, Faculty of Geosciences, Utrecht University, Utrecht, 3584 CB, The Netherlands

11 *Correspondence to:* Madleen Grohganz (madleen.grohganz@nioz.nl)

12 **Abstract.** Tipping points are thresholds at which a system, often abruptly and irreversibly, transitions from a stable
13 state to a contrasting one. Crossing such critical boundaries poses a risk to Earth system stability and may have
14 catastrophic consequences. This is especially relevant, as current climate change is destabilizing Earth subsystems,
15 potentially bringing them closer to tipping points. Thus, it is important to be able to detect approaching tipping points
16 in the Earth's system, which can be achieved through calibration on palaeo-records. Recently, new deep learning (DL)
17 methods have been established that are able to confidently and quantitatively identify different types of critical
18 transitions characterised by their abruptness and (ir)reversibility. Based on this, we develop a new (simplified) DL
19 classifier focusing on the quantitative detection of catastrophic tipping points (fold bifurcations) in the Earth system.
20 Our approach reduces computational demand and improves performance, especially for short timeseries. We first test
21 the new classifier's performance on synthetic data and subsequently on different existing Cenozoic proxy records. Our
22 DL results are compared to the results from previous studies applying generic early warning signals (EWS), which can
23 detect approaching transitions qualitatively but cannot distinguish bifurcation types (abruptness and (ir)reversibility of
24 the transition). Our DL classifier enables us to identify how abrupt and (ir)reversible an approaching transition is,
25 which is important for tipping point risk assessment and mitigation. Results are generally consistent between generic
26 EWS from previous studies and our DL approach and fit with what is known from the geological context. We note that
27 some results are dependent on the length of the classifier used and the time interval investigated before the bifurcation.
28 We implement an out-of-distribution (OOD) detection method to reduce the misclassification of non-catastrophic
29 bifurcations as catastrophic tipping points. Combined with the binary DL classifier, this approach enables reliable,
30 quantitative detection of catastrophic tipping points in Earth system records.

31

32 **Keywords:** tipping point, nonlinear dynamics, Earth system, machine learning, early warning signals

33

34



35 **1 Introduction**

36 System Earth is vulnerable to tipping points, critical thresholds at which gradual changes in external forcings
37 destabilize key components, pushing them towards often abrupt and irreversible transitions (Scheffer et al., 2009; Van
38 Nes et al., 2016; Lenton et al., 2008). Prominent examples include the Atlantic Meridional Overturning Circulation
39 (AMOC) (Rahmstorf, 1995; Hawkins et al., 2011; Van Westen and Dijkstra, 2023; Van Westen et al., 2024), the
40 Greenland and Antarctic ice sheets (Van Breedam et al., 2023; Robinson et al., 2012; Garbe et al., 2020), and other
41 Earth subsystems that exhibit catastrophic tipping behavior under sustained anthropogenic climate forcing (Lenton et
42 al., 2025; Lenton et al., 2008). Crossing thresholds in these systems can trigger drastic global impacts affecting millions
43 of people worldwide, such as abrupt shifts in temperature and precipitation regimes, accelerated sea level rise
44 (Drijfhout, 2015; Jackson et al., 2015; Armstrong McKay et al., 2022), and even tipping cascades, i.e. a domino-like
45 collapse of several interconnected Earth systems (Wunderling et al., 2021; Kriegler et al., 2009).

46 Critical transitions are defined by a loss of stability. Changing external forcings, such as greenhouse gas increases or
47 ocean temperature changes, destabilize the system and push it towards a critical boundary, a so-called bifurcation point.
48 The resulting shift may unfold gradually or discontinuously, with the latter defining a catastrophic tipping point: rapid,
49 self-perpetuating, and often irreversible due to hysteresis (Milkoreit et al., 2018; Scheffer et al., 2009; Lenton et al.,
50 2008). Hysteresis means that, even if forcings are removed or reversed to the previous level after passing a tipping
51 point, the system remains locked in its new state until conditions are reduced even further (overcompensated), thereby
52 complicating recovery efforts (Crawford, 1991).

53 Current approaches to tipping point detection rely on so-called early warning signals (EWS) (Scheffer et al., 2009).
54 EWS are based on the observation that a system approaching a tipping point loses resilience and becomes slower at
55 recovering from small perturbations, a behaviour called critical slowing down (CSD) (Wissel, 1984; Scheffer et al.,
56 2009). Different statistical techniques can be used to infer CSD from timeseries leading up to tipping points, e.g.
57 increasing lag-1 autocorrelation and variance. Using these CSD-based EWS, several studies have identified
58 approaching critical transitions in palaeoclimatic timeseries as well as observational climate timeseries, e.g. related to
59 marine (de)oxygenation (Hennekam et al., 2020a), the Greenland ice sheet (Boers and Rypdal, 2021), the AMOC
60 (Boers, 2021), Dansgaard-Oeschger events (Boers, 2018) and major Cenozoic climate transitions (Boettner et al.,
61 2021). However, these generic EWS are limited to qualitatively flagging approaching critical transitions in general.
62 They cannot distinguish between smooth, easily reversible and discontinuous, difficult to reverse transitions. However,
63 this aspect is particularly important for Earth Systems, where reversibility of transitions have major implications. To
64 better assess the potential impacts of a critical transition on the Earth system, we need to move away from the
65 qualitative nature of generic EWS to a more quantitative approach, that can confidently detect catastrophic tipping
66 points of an abrupt and irreversible nature.

67 Recent studies have applied machine learning to quantitatively forecast tipping points. These studies trained deep
68 learning classifiers to recognize time series leading up to different types of critical transitions/bifurcations (Bury et al.,
69 2021; Bury et al., 2023; Deb et al., 2022). This approach outperformed generic EWS, i.e. traditional statistical



70 indicators (lag-1 autocorrelation and variance) in signaling an upcoming tipping point. Importantly, DL classifiers
71 provide the ability to specifically identify systems with rapid and irreversible transitions, which conventional CSD-
72 based indicators alone cannot achieve. The most critical bifurcation type identified in Earth systems is the fold
73 bifurcation, which is associated with abrupt, irreversible tipping points (Scheffer et al., 2009; Lenton et al., 2025;
74 Lenton et al., 2008). There are other bifurcation types that are also subject to critical slowing down (Kéfi et al., 2014),
75 but involve a smooth transition, like the transcritical bifurcation. Transcritical bifurcations are non-generic and
76 generally represented by simple models with strict bounds and are relatively rare in nature (an exception being e.g. the
77 onset of epidemics) (Arnold, 2012; Seydel, 2009; Kermack and Mckendrick, 1927). But they are particularly difficult
78 to distinguish from fold bifurcations due to their mathematical similarity (both have quadratic, one-dimensional normal
79 forms (Kuznetsov, 1998)) and thus are regularly misclassified as fold bifurcations by current DL algorithms. We want
80 to limit the amount of misclassified transcritical bifurcations and ensure, that the classifier confidently detects
81 bifurcations that are abrupt and irreversible (fold bifurcation).

82 Here, we build on existing multi-class DL classifiers (Bury et al., 2023; Bury et al., 2021) and train a new binary
83 classifier which focuses exclusively on the detection of fold bifurcations (abrupt, catastrophic tipping points). Multi-
84 class classifiers are trained on multiple different types of bifurcations, while binary classifiers specifically focus on
85 two bifurcation types of interest. We opted for a binary classifier over a multi-class classifier for two reasons: (i) it is
86 computationally less expensive to train, and (ii) we expect it to have better performance on the specific task of detecting
87 a fold bifurcation. Additionally, we implement out-of-distribution (OOD) detection to flag misclassified transcritical
88 bifurcations, enhancing reliability. We test the performance of our binary classifier with OOD on synthetic data and
89 compare its performance metrics to the multi-classifier approach (Bury et al., 2021). We then apply our approach to
90 different Earth system proxy records, that represent instances of past transitions in temperature, carbon cycling and
91 marine (de)oxygenation. These records reflect processes acting on different timescales, from several thousand to
92 millions of years and have previously been shown to exhibit generic CSD-based EWS (lag-1 autocorrelation, variance)
93 (Setty et al., 2023; Boettner et al., 2021; Hennekam et al., 2020a). Applying our new DL classifier and OOD detection
94 we assess abrupt tipping behaviour associated with catastrophic fold bifurcations for these records in order to evaluate
95 the potential added value of the approach proposed here.

96

97 **2 Material and methods**

98 **2.1 Training data for the DL classifier**

99 We trained two types of binary classifiers suitable to analyse timeseries of up to 500 datapoints (length 500 classifier)
100 and timeseries of up to 1500 datapoints (length 1500 classifier). For training the DL classifiers we used the synthetically
101 generated training datasets of length 500 and length 1500 respectively from Bury et al. (2021). The length 500 training
102 dataset consists of 125,000 null time series and fold time series each, the length 1500 training dataset of 50,000 null
103 and fold timeseries each. These timeseries are derived from simulations of randomly generated, two dimensional
104 systems, that exhibit local bifurcations. Null timeseries are obtained from simulations with fixed parameters, fold



105 bifurcations from simulation with varying parameter towards the bifurcation. To obtain the residual time series, they
106 were detrended using Lowess smoothing with a span of 0.2 (smoothing parameter indicating the width of the moving
107 window used to calculate the residuals; in our case 20% of the original timeseries datapoints). Each residual time series
108 was normalized by dividing each time series data point by the average absolute value of the residuals across the entire
109 time series in order to remove scale dependence and allow comparison across different timeseries

110 **2.2 DL algorithm architecture and training**

111 For our DL classifier training we adapted the CNN-LSTM DL algorithm from Bury et al. (2021) (see their paper for
112 more details) with one convolutional (CNN) layer and two recurrent (LSTM) layers. The algorithm was trained for
113 500 epochs with a learning rate of 0.0005. Due to the high number of synthetically generated timeseries, we used a
114 train/validation/test split of 0.95/0.04/0.01. Two variants (model 1 and model 2) of the length 500/length 1500
115 algorithm were trained on censored versions of the training time series. Model 1 was trained on time series, that were
116 censored by padding them with a variable number of zeros at both the left and right (beginning and end of the time
117 series). This results in time series of various lengths (as short as 50), that can represent any part of the time series; not
118 necessarily the time phase before the transition, but also the middle or beginning sections. Model 2 was trained on a
119 censored version of the training set, where the time series were padded with zeros only on the left (at the beginning of
120 the timeseries). This results in time series of various lengths (as short as 50); but all of them directly leading up to the
121 bifurcation. Due to their different censoring, model 1 excels at picking up early warning signals from the early parts
122 of the timeseries and model 2 rather at detecting the specific type of bifurcation close to the actual transition. To
123 leverage and integrate the capabilities of both models, we take the average prediction of five model 1 and five model
124 2 classifiers at each point to generate the reported results (ensemble) for each classifier (length 500 and length 1500).
125 Prior to application on empirical systems we evaluate the performance of our binary classifier on the test portion of
126 the synthetically generated dataset (see results section).

127 For each of the length 500 and length 1500 classifiers we calculated different performance metrics: precision (how
128 many positive predictions are true positives), recall (how many of the true positives are detected), and f1 score (a
129 combined measure of precision and recall). The final performance values are the means of these metrics for an
130 ensemble of five model 1 and five model 2 classifiers for length 500 and length 1500 respectively.

131 **2.3 Empirical systems for application of the binary DL classifier**

132 We then apply the binary DL classifier to three different empirical systems from the Cenozoic and compare its results
133 with previous studies focussing on generic EWS (lag-1 autocorrelation and variance):

- 134 1) The stable isotope record of the Paleocene and Early Eocene, in which several palaeoclimatic events
135 associated with abrupt perturbations of the carbon cycle ($\delta^{13}\text{C}$) and changes in temperature ($\delta^{18}\text{O}$) can be
136 identified (PETM, ETM2 and ETM3 events) (Fig. 1). The PETM is the most intense warming event
137 (hyperthermal) of the Cenozoic accompanied by a rapid release of isotopically light carbon into the
138 atmosphere-ocean system, global warming and significant ocean acidification; it was followed by the similar,
139 but less intense hyperthermal events ETM2 and ETM3 (Westerhold et al., 2020; Zachos et al., 2005; Sexton

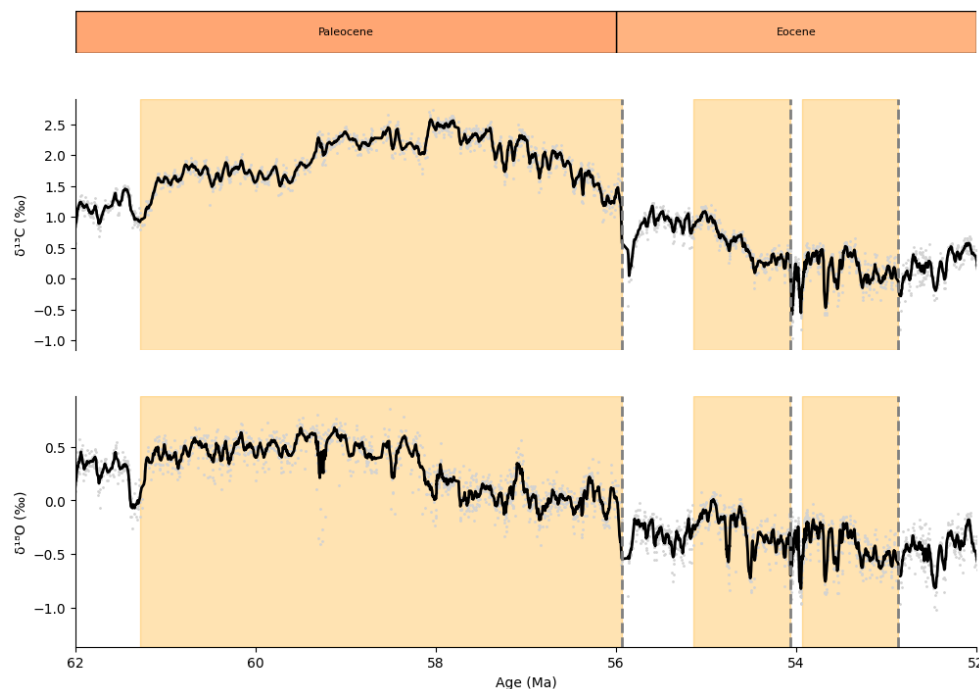


140 et al., 2011; Lauretano et al., 2018b). To assess the abruptness of these transitions we combine the isotope
 141 records from ODP Site 208-1262 (Barnet et al., 2019; Barnet et al., 2018) (PETM and ETM2) and ODP Site
 142 208-1263 (Lauretano et al., 2018b; Lauretano et al., 2018a) (ETM3). As our DL classifier requires a time
 143 series with equally spaced data points, the $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ data are interpolated to regular sampling steps (using
 144 the nearest neighbor method). We apply the same analysis timeframes (see Table 1) and data pre-processing
 145 methods as Setty et al. (2023) and compare the results of our analyses to the generic EWS from their study.
 146 For timeseries longer than 500 datapoints (PETM) we apply both the length 1500 classifier (on the complete
 147 timeseries) and the length 500 classifier (on the last 500 points before the bifurcation).

148 **Table 1: Beginning and end dates of analysed events (based on Setty et al., 2023) and number of data points in**
 149 **the respective time series.**

Event	Transition [Ma]	Analysis timeframe [Ma]	Number of data points (and applied classifier)
Paleocene/Eocene Thermal Maximum (PETM)	55.935	61.28–55.935	1397 (length 1500 classifier) and 500 (length 500 classifier)
Eocene Thermal Maximum 2 (ETM2)	54.06	55.141–54.06	416 (length 500 classifier)
Eocene Thermal Maximum 3 (ETM3)	52.87	53.936–52.87	411 (length 500 classifier)

150



151

152 **Figure 1: $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ isotope record of the Paleocene and Early Eocene with PETM, ETM2 and ETM3 events**
 153 **marked with dashed lines, areas shaded in yellow indicate the analysed time intervals prior to the events.**

154

155 2) The CENOGRID dataset, a high-resolution $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ isotope record covering the Cenozoic (Westerhold
 156 et al., 2020). The CENOGRID dataset is a continuous composite of benthic foraminifer isotope records,
 157 derived from 14 Atlantic and Pacific ocean drilling cores mainly from low to mid latitudes (Westerhold, 2020).
 158 Several Cenozoic palaeoclimatic events associated with abrupt changes in temperature ($\delta^{18}\text{O}$) and abrupt
 159 perturbations of the carbon cycle ($\delta^{13}\text{C}$) can be identified in this record. Here, we focus on the
 160 Eocene/Oligocene Transition, Oligocene/Miocene Transition and Middle Miocene Climate Transition, events
 161 associated with significant positive $\delta^{18}\text{O}$ isotope excursions indicating shifts from a warm greenhouse climate
 162 towards a colder icehouse climate (Fig. 2). These transitions are related to different steps in the build-up and
 163 expansion of the Antarctic ice sheet (Lear et al., 2008; Wilson et al., 2008; Zachos et al., 2001; Shevenell et
 164 al., 2004), potentially representing tipping points of the ice sheet. To assess the abruptness of these events we
 165 use the binned and interpolated isotope records (Westerhold (2020), S34; see also SI of the original paper
 166 Westerhold et al. (2020)). We adapt the analysis timeframes (see Table 2) from Boettner et al. (2021), which
 167 are optimized to cover the highest data point density leading up to the transition. The results of our analyses
 168 are compared to the generic EWS from Boettner et al. (2021), who use non-equidistant time series from the
 169 binned, but non-interpolated version of the CENOGRID record. For timeseries longer than 500 datapoints

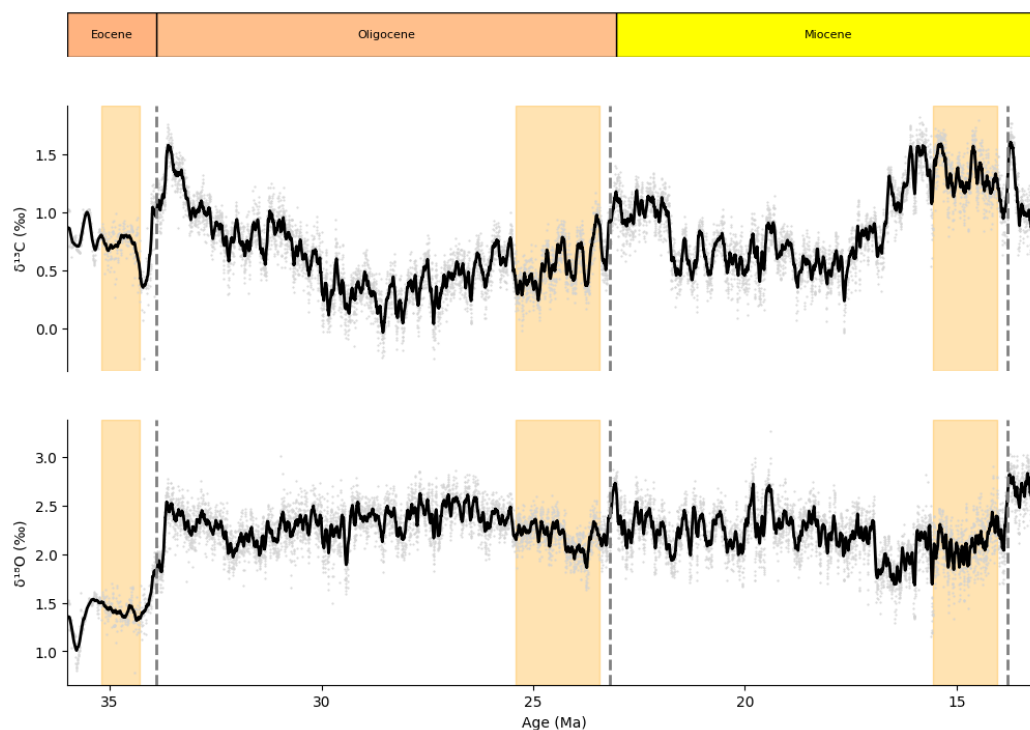


170 (Oligocene/Miocene Transition and Middle Miocene Transition) we apply both the length 1500 classifier (on
 171 the complete timeseries) and the length 500 classifier (on the last 500 points before the bifurcation).

172 **Table 2: Beginning and end dates of analysed events (based on Boettner et al., 2021) and number of data points**
 173 **in the respective time series.**

Event	Transition [Ma]	Analysis timeframe [Ma]	Number of data points (and applied classifier)
Eocene/Oligocene Transition	33.9	35.2-34.305	180 (length 500 classifier)
Oligocene/Miocene Transition	23.2	25.424-23.44	993 (length 1500 classifier) and 500 (length 500 classifier)
Middle Miocene Climate Transition	13.8	15.56-14.04	761 (length 1500 classifier) and 500 (length 500 classifier)

174



175



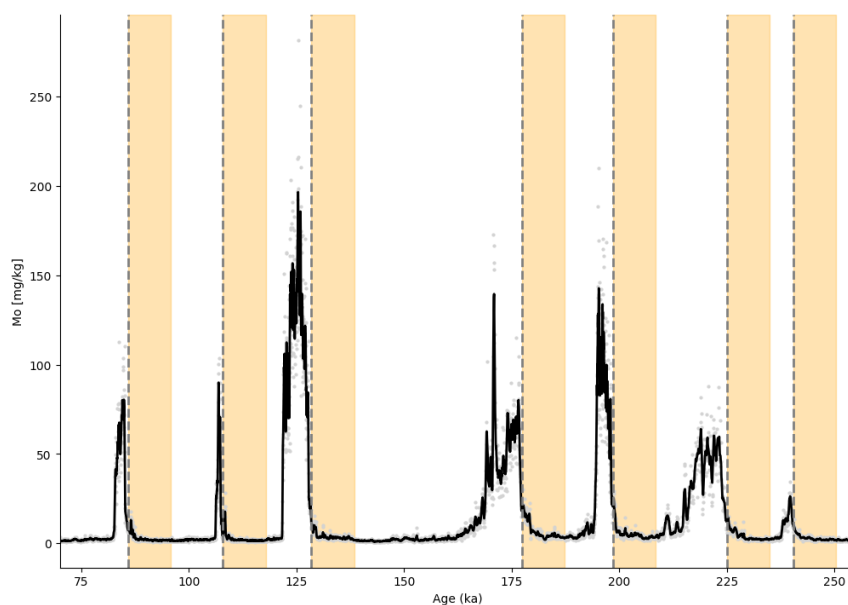
176 **Figure 2: CENOGRID $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ isotope record with Eocene/Oligocene Transition, Oligocene/Miocene**
 177 **Transition and Middle Miocene Climate Transition events marked with dashed lines, areas shaded in yellow**
 178 **indicate the analysed time intervals prior to the events.**

179 3) The eastern Mediterranean sapropel record. Sapropels are organic-rich sediment layers occurring regularly in
 180 eastern Mediterranean sediments and represent changes from an oxic to anoxic state (and back) as a result of
 181 relatively gradual increasing freshwater inputs (Hennekam et al., 2020a). To assess the abruptness of these
 182 oxic-to-anoxic transitions we use the anoxia proxy Mo (calibrated XRF-scanning data) from core 64PE406E1
 183 (Hennekam et al., 2020b) (Fig. 3). We use the original dataset provided; interpolation is not required, as the
 184 spacing of datapoints is almost regular. We adapt the following analysis timeframes (see Table 3) and compare
 185 our results to the generic EWS from Hennekam et al. (2020b). As all sapropel timeseries are shorter than 500
 186 datapoints we exclusively apply the length 500 classifier.

187 **Table 3: Beginning and end dates of analysed sapropel events (based on Hennekam et al., 2020b) and number**
 188 **of data points in the respective time series.**

Event	Transition [ka]	Analysis timeframe [ka]	Number of data points (and applied classifier)
S9	240.378	250.318-240.378	201 (length 500 classifier)
S8	225.039	234.923-225.039	232 (length 500 classifier)
S7	198.533	208.494-198.533	289 (length 500 classifier)
S6	177.308	187.225-177.308	211 (length 500 classifier)
S5	128.393	138.361-128.393	324 (length 500 classifier)
S4	107.823	117.818-107.823	435 (length 500 classifier)
S3	85.826	95.823-85.826	376 (length 500 classifier)

189



190

191 **Figure 3: Eastern Mediterranean sapropel record of the anoxia proxy Mo with analysed S3 to S9 sapropel events**
192 **marked with dashed lines, areas shaded in yellow indicate the analysed time intervals prior to the events.**

193 For analysis of the results we adapted a version of ewstools (Bury, 2023; Bury et al., 2021) for our newly developed
194 binary classifier. The script takes timeseries as inputs, which are detrended with the same filter (Lowess filter of span
195 0.2) as the synthetic timeseries the classifier was trained on (Dablander and Bury, 2022). We then apply our binary
196 classifiers using the `apply_classifier_inc()` function defined in ewstools with an `inc` value of 10. This calculates DL
197 probabilities for the input timeseries as they are revealed stepwise in time intervals of 10 until they reach the -critical
198 transition. Results are plotted as DL probability over time up to the transition. Individual DL probability graphs are
199 calculated for the five classifiers of model 1 and the five classifiers of model 2. Ensembled results are plotted as the
200 average DL probabilities of all models at a specific timestep. A fold or null transition is assigned if the DL probability
201 output of the ensembled models is >0.5 respectively. We visualise these results as heatmaps, in which we plot the DL
202 probabilities of the dominant transition type (fold or null) at the last timestep before the transition. Results are shown
203 for variables ($\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ isotope values, Mo concentrations) across the investigated events.

204 **2.4 OOD detection methods**

205 To evaluate how confident our classifier is in detecting catastrophic fold bifurcations and how many transcritical
206 bifurcations are misclassified as fold, we apply OOD detection methods. We use two approaches, (i) the classical
207 softmax confidence score approach (Hendrycks and Gimpel, 2016) and (ii) a new energy-based approach (Liu et al.,
208 2020).

209 The softmax confidence scores are the same as the aforementioned DL probability outputs, that are directly derived
210 from the application of the classifier. These values represent the model's prediction confidence across the fold and null



211 classes for an individual timeseries. We ran tests with the softmax probability to investigate how overconfident our
212 classifier is in its results. We applied our ensembled length 500 binary classifiers (five model 1 and five model 2
213 classifiers) on 100 synthetically generated transcritical timeseries from the test dataset (see Bury et al., 2021). Our
214 classifiers were only trained on null and fold bifurcations and had no previous contact with transcritical timeseries
215 during training. We investigated how many of these transcritical timeseries get erroneously classified as fold
216 bifurcations by our binary classifier.

217 We also applied an energy-based out-of-distribution (OOD) detection approach (Liu et al., 2020) as a posthoc method
218 to distinguish if timeseries, that were originally classified as fold by our model are true fold bifurcations or
219 misclassified transcritical bifurcations. Energy-based OOD detection has been shown to be less susceptible to
220 overconfidence than the softmax confidence score approach and better at identifying OOD samples in image detection
221 applications (Liu et al., 2020). We adapted this energy-based OOD detection method originally developed for image
222 networks and implement it for timeseries. Energy-based OOD detection relies on energy scores calculated from an
223 energy function. The energy function takes the raw logits outputs from a classifier and combines them into a scalar,
224 that measures how strongly the model recognizes the input timeseries as belonging to any known class (the classifier
225 was trained on). The higher the calculated energy score, the more confident the classifier is about this timeseries being
226 in-distribution (ID) (representing a bifurcation class known to our classifier, in this case fold bifurcation). The lower
227 the calculated energy score, the more likely it is out-of-distribution (OOD) (representing a bifurcation class unknown
228 to our classifier, in this case transcritical). The resulting energy scores are plotted in an energy density histogram,
229 showing the distribution of energy score values of ID and OOD timeseries (see Fig. 5).

230 We calculate such energy density histograms for our length 500 and length 1500 binary classifiers. We use an ensemble
231 of five binary classifiers (model 2) for length 500 as well as length 1500. Our ensembles only contain model 2, as here
232 we focus on determining the type of bifurcation (fold vs transcritical) and not generally detecting early warning (which
233 model 1 works best for). In addition, the empirical data we later apply the OOD detection to, is also in the form of the
234 input data for model 2, i.e. timeseries going right up to the bifurcation. As input for our energy density histograms we
235 take fold and transcritical timeseries from the test dataset of length 500 (1250 timeseries each) and length 1500 (500
236 timeseries each) (see Bury et al., 2021). For these test dataset timeseries their bifurcation type is known, so the resulting
237 energy density histograms represent the energy score distributions of ID (fold) and OOD (transcritical) timeseries for
238 the length 500 and length 1500 classifier respectively (see Fig. 5). An energy threshold is set, that divides ID and OOD
239 energy scores. We chose a threshold value of 2.9 (length 500) and 3.85 (length 1500) respectively (yellow lines in Fig.
240 5), as the number of correctly identified ID and OOD timeseries was equal for these energy scores.

241 We then compare energy scores of our empirical timeseries (identified as fold by our binary classifier) to our energy
242 density histograms. For empirical timeseries shorter than 500 points and 1500 datapoints respectively, the timeseries
243 were padded with zeros at the front. All empirical timeseries were detrended using a Lowess filter of span 0.2 before
244 analysis to match the preprocessing of the test dataset timeseries (used as a basis for the energy density histograms). If
245 the energy score values calculated for the empirical timeseries fall above the energy threshold, i.e. within the ID (fold)



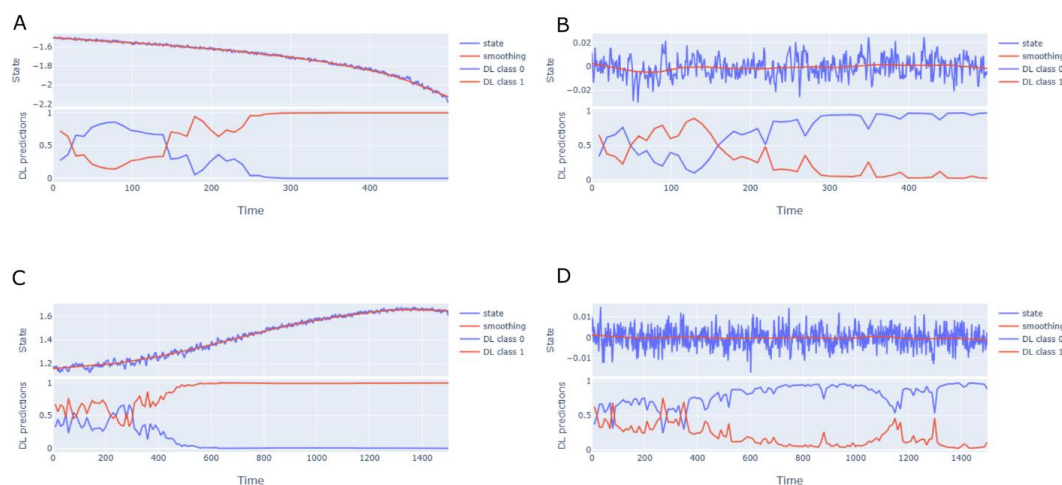
246 part of the energy density histogram, they constitute true fold bifurcations; otherwise they are flagged as misclassified
247 transcritical bifurcations.

248 To test for noise as the underlying cause of the bimodal nature of the ID and OOD energy density histograms (see Fig.
249 5), we calculated the standard deviation of the underlying test dataset timeseries. We applied Kruskal Wallis tests to
250 determine statistically significant differences in the standard deviations (noise levels) of ID and OOD timeseries above
251 and below the energy score threshold.

252 3 Results

253 3.1 Performance of the binary DL classifier on the synthetic test data

254 We evaluated the performance of our binary classifiers of length 500 (Fig. 4A and B) and length 1500 (Fig. 4C and D).
255 The length 500 classifier achieved 86.3% precision, 85.6% recall, and an F1 score of 85.5% (Table 4). The length 1500
256 classifier performed slightly better (87.8% precision, 87.5% recall, F1 score 87.4%, see Table 4). Both our binary
257 length 500 and length 1500 classifiers are able to correctly classify most timeseries, with a relatively low and balanced
258 amount of false negatives and false positives. We compare the performance metrics of our binary classifiers with multi-
259 class classifiers (Bury et al., 2021) (see Table 4). For shorter timeseries (length 500), our binary classifier outperformed
260 the multiclass alternative. However, for length 1500 timeseries the multiclass classifier achieved slightly better
261 performance values.



262

263 **Figure 4: Examples of analyses with the binary classifier ensemble of length 500 on a synthetic fold timeseries**
264 **(A) and synthetic null timeseries (B), and with the binary classifier ensemble of length 1500 on a synthetic fold**
265 **timeseries (C) and a synthetic null timeseries (D); with state as the original timeseries, DL class 0 as average of**
266 **ensembled model results for null bifurcation and DL class 1 as average of ensembled model results for fold**
267 **bifurcation.**



268 **Table 4: Performance metrics of our binary classifier (length 500 and length 1500) compared to performance**
 269 **metrics multiclass classifiers (length 500 and length 1500) (Bury et al., 2021)**

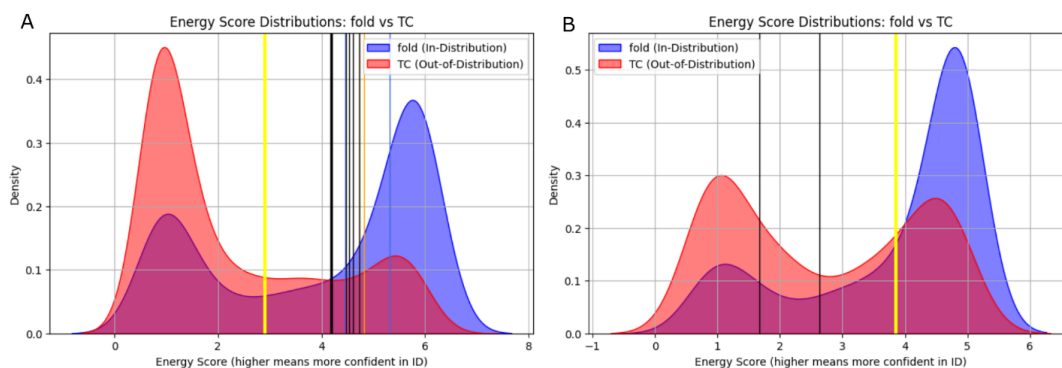
	Binary classifier (length 500)	Binary classifier (length 1500)	Multiclass classifier (length 500) (Bury et al., 2021)	Multiclass classifier (length 1500) (Bury et al., 2021)
Precision	86.3%	87.8%	84.4%	88.3%
Recall	85.6%	87.5%	84.2%	88.3%
F1 score	85.5%	87.4%	84.2%	88.2%

270

271 **3.2 OOD detection on synthetically generated test data**

272 We test different types of OOD detection methods in order to identify transcritical bifurcations, that are misclassified
 273 as fold bifurcations. The softmax confidence score method reveals that 76% of synthetically generated transcritical
 274 timeseries are misclassified as fold by our binary classifier (see supplementary data). To account for this high number
 275 of misclassification we apply an energy-based OOD detection approach. We calculate energy density histograms for
 276 our length 500 and length 1500 binary classifier (Fig. 5) (for more details see the material and methods section).
 277 Energy-based OOD detection allows for a clear distinction between the energy scores of ID (fold) and OOD
 278 (transcritical) timeseries for length 500 (Fig. 5A), and to a lesser degree also for length 1500 (Fig. 5B). This leads to
 279 about 68% (length 500) and 66% (length 1500) of fold and transcritical timeseries being correctly identified based on
 280 their energy scores, thus reducing misclassification to 32% (length 500) and 34% (length 1500) with our energy-based
 281 OOD detection method. However, there remains a certain overlap between ID and OOD in the energy density
 282 histogram due to the bimodal nature of the curves, more pronounced in length 1500 (Fig. 5B) than length 500 (Fig.
 283 5A).

284



285 **Figure 5: Energy density histograms (A) based on the length 500 test dataset of synthetically generated**
 286 **transcritical (TC) timeseries (OOD) and fold timeseries (ID), the yellow line indicates the threshold between**
 287 **OOD and ID, black lines mark energy scores of empirical sapropel timeseries, blue lines empirical ETM2**



288 **timeseries and orange lines empirical Oligocene/Miocene timeseries (see below); (B) based on the length 1500**
 289 **test dataset of synthetically generated timeseries, the yellow line indicates the threshold between OOD and ID,**
 290 **black lines mark energy scores of empirical PETM timeseries (see below).**

291 We tested for noise as the underlying cause of this bimodality, using standard deviation as a measure of noise. For both
 292 length 500 and length 1500, timeseries with higher energy scores (above the energy score threshold) are associated
 293 with higher median noise values compared to the ones with lower energy scores (below the energy threshold) (Table
 294 5). High energy scores and thus confidence in identification generally seem to be associated with timeseries with higher
 295 noise for both ID (fold) and OOD (transcritical) timeseries.

296 **Table 5: Mean noise levels of ID and OOD timeseries of different lengths above and below the energy threshold**
 297 **and results of Kruskal-Wallis test to determine significance of noise differences between the groups.**

	Mean noise (std) above energy threshold	Mean noise (std) below energy threshold	p-value Kruskal-Wallis test
ID (length 500)	0.014	0.011	p-value = 5.6×10^{-25}
OOD (length 500)	0.016	0.012	p-value = 3.6×10^{-17}
ID (length 1500)	0.015	0.011	p-value = 3.6×10^{-09}
OOD (length 1500)	0.018	0.013	p-value = 7.5×10^{-12}

298

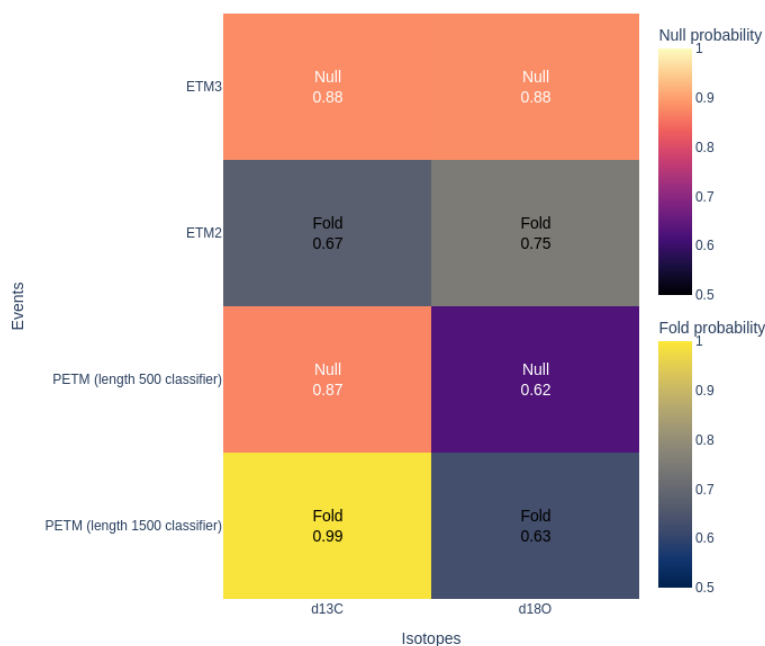
299 **3.3 Application of the binary classifier to Cenozoic proxy record**

300 **3.3.1 Paleocene/Eocene isotope record (PETM, ETM2, ETM3)**

301 The PETM, ETM2 and ETM3 events all show pronounced negative $\delta^{13}\text{C}$ as well as $\delta^{18}\text{O}$ excursions indicating strong
 302 carbon cycle perturbations and warming (see Fig. 1). The PETM is considered the strongest of the three events,
 303 followed by ETM2 and ETM3. We applied our binary DL classifier to these Paleocene/Eocene isotope events. Fig. 6
 304 summarises the DL predictions and probabilities for the $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ isotope records of the PETM, ETM2 and ETM3
 305 events. For the PETM our length 1500 classifier predicts fold bifurcations in both the $\delta^{13}\text{C}$ (0.99) and $\delta^{18}\text{O}$ (0.63)
 306 isotope records, indicating an approaching tipping point in the carbon cycle, as well as in temperature (albeit with a
 307 weaker signal). This supports the view of the PETM as an event mainly driven by a catastrophic tipping point in the
 308 carbon cycle. However, applying our length 500 binary classifier returns null predictions in both isotope records (0.87
 309 for $\delta^{13}\text{C}$ and 0.62 for $\delta^{18}\text{O}$), indicating that there is no fold bifurcation approaching. For the ETM2 (analysed with the
 310 length 500 classifier) our classifier predicts fold bifurcations for both isotope records (0.67 for $\delta^{13}\text{C}$ and 0.75 for $\delta^{18}\text{O}$),
 311 indicating approaching tipping points. Our results support the ETM2 as a tipping point event driven by rapid changes
 312 in both the carbon cycle and temperature. The ETM3 event (analysed with the length 500 classifier) returns null
 313 predictions (0.88 for both $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$), both with high DL probabilities, indicating no approaching bifurcations.
 314 This corroborates ETM3 being the weakest of the Early Eocene events, which didn't go through a catastrophic tipping
 315 point.



316



317

318 **Figure 6: Heatmap of DL classifier predictions and probabilities for $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ isotope records of the PETM,**
 319 **ETM2 and ETM3 events, including results for different classifier lengths.**

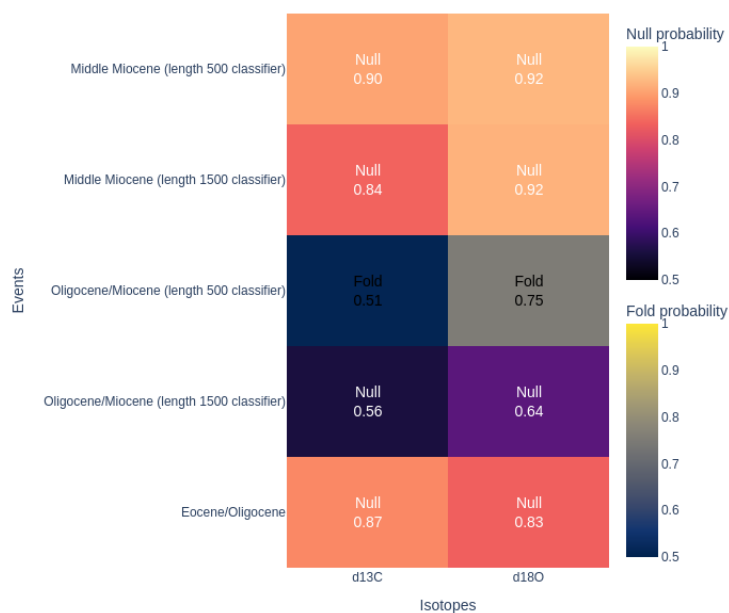
320

321 **3.3.2 CENOGRID isotope record**

322 The Eocene/Oligocene Transition, the Oligocene/Miocene Transition and the Middle Miocene Transition all show
 323 positive $\delta^{13}\text{C}$ as well as $\delta^{18}\text{O}$ excursions indicating general cooling and ice volume increase (see Fig. 2). Fig. 7
 324 summarises our classifier's DL predictions and probabilities for the $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ isotope records of the analysed
 325 events. For the Eocene/Oligocene Transition the 500 length classifier returns null predictions for both the $\delta^{13}\text{C}$ (0.87)
 326 and $\delta^{18}\text{O}$ records (0.83), with high DL probabilities. The Eocene/Oligocene Transition is considered the first onset of
 327 the Antarctic glaciation, however no tipping point was crossed, the system did not yet transition directly into a fully
 328 glacial state (Galeotti et al., 2016). For the Oligocene/Miocene Transition our 500 length classifier returns fold
 329 predictions for both isotope records (0.75 for $\delta^{18}\text{O}$ and 0.51 for $\delta^{13}\text{C}$), pointing to substantial perturbations of the
 330 temperature record (and to a lesser degree also the carbon cycle). We link this to a tipping point being crossed in the
 331 Antarctic ice sheet subsystem, locking the system into a cool state with a stable Antarctic icesheet. However, when



332 applying our 1500 length classifier, null predictions are returned, albeit with low DL probabilities (0.56 for $\delta^{13}\text{C}$ and
 333 0.64 for $\delta^{18}\text{O}$). For the Middle Miocene Transition both our 500 length and 1500 length classifier return null predictions
 334 for both isotope records, all with high DL probabilities (500 length classifier: 0.90 for $\delta^{13}\text{C}$ and 0.92 for $\delta^{18}\text{O}$, 1500
 335 length classifier: 0.84 for $\delta^{13}\text{C}$ and 0.92 for $\delta^{18}\text{O}$), indicating no tipping points associated with the Middle Miocene
 336 records. During the Middle Miocene Climate Transition temperature continued to decrease, the ice sheets stabilized
 337 and the system got fully locked in an icehouse mode, but without crossing another tipping point.



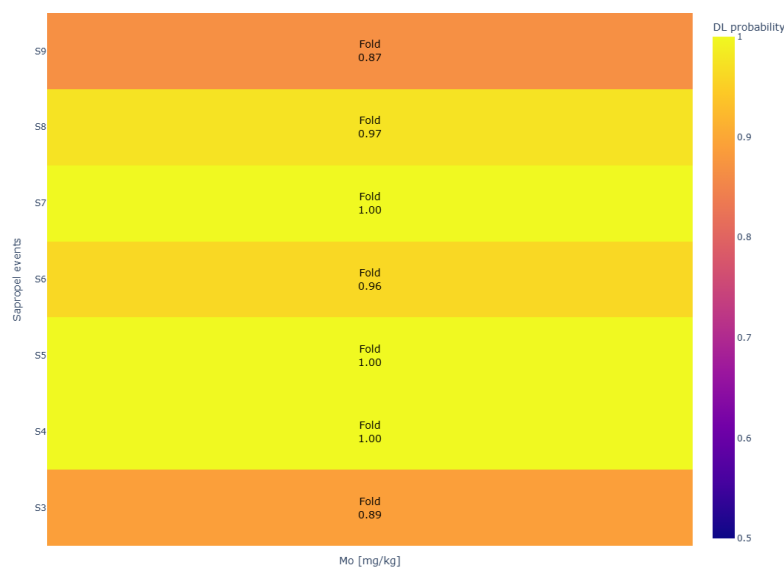
338

339 **Figure 7: Heatmap of DL classifier predictions and probabilities for $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ isotope records of the**
 340 **Eocene/Oligocene Transition, Oligocene/Miocene Transition and Middle Miocene Climate Transition, including**
 341 **results for different classifier lengths.**

342

343 3.3.3 Eastern Mediterranean sapropel data

344 Our binary fold classifier was applied to the Eastern Mediterranean sapropel record of S3 to S9 from core 64PE406E1.
 345 Fig. 8 summarises the classifier's DL predictions and probabilities for the Mo proxy records of the analysed events.
 346 For all sapropel timeseries (analysed with the length 500 classifier) fold bifurcations are predicted. DL probabilities
 347 are generally high (>0.95), except for the weakest sapropel events S3 (0.89) and S9 (0.87). This clearly indicates a
 348 tipping point from oxic to anoxic conditions associated with all investigated Eastern Mediterranean sapropel layers.



349

350 **Figure 8: Heatmap of DL classifier predictions and probabilities for the Mo anoxia proxy record of various**
351 **Eastern Mediterranean sapropel events.**

352

353 3.6 OOD detection in empirical timeseries data

354 We subsequently applied our energy based OOD detection approach to the empirical timeseries identified as fold
355 bifurcations by our binary classifier to flag any potentially misclassified transcritical timeseries. We evaluate empirical
356 timeseries with length 500 or shorter, i.e. the sapropel timeseries, ETM2 timeseries ($\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ isotope record) and
357 Oligocene/Miocene timeseries ($\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ isotope record), against the length 500 energy density histogram. The
358 calculated energy scores of all the investigated empirical timeseries fall above the energy score threshold and thus
359 within the ID part of the energy score distribution (black, orange and blue vertical lines in Fig. 5A). The sapropel Mo
360 timeseries (S3 to S9), the ETM2 timeseries ($\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ isotope record) and the Oligocene/Miocene timeseries ($\delta^{13}\text{C}$
361 and $\delta^{18}\text{O}$ isotope record) were all correctly identified by our binary classifier as genuine fold bifurcations. We evaluated
362 longer empirical timeseries of up to 1500 data points, i.e. the PETM timeseries ($\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ isotope record), against
363 the length 1500 based energy density histogram. The calculated energy scores for both the PETM $\delta^{18}\text{O}$ and $\delta^{13}\text{C}$
364 timeseries fall below the energy threshold, flagging both PETM timeseries as misclassified transcritical bifurcations
365 instead of a genuine fold bifurcation.

366



367 4 Discussion

368 4.1 Performance of our binary classifier with OOD in comparison to other EWS methods

369 We introduce a new binary classifier focused on the detection of fold bifurcations in Earth systems. Fold bifurcations
370 are the most dominant transitions in these systems and are widely associated with Earth system tipping points (e.g. in
371 the AMOC (Rahmstorf, 1995; Van Westen et al., 2025; Stommel, 1961) or the polar ice sheets (Robinson et al., 2012;
372 Garbe et al., 2020)). Our new classifier is able to confidently identify fold bifurcations in Earth system applications,
373 as shown by our empirical test cases. This constitutes a significant improvement over generic EWS (autocorrelation
374 and variance based), as it allows to quantitatively detect approaching critical thresholds and identify their bifurcation
375 type instead of providing only qualitative indication of an otherwise unspecified approaching transition.

376 Compared to other already existing multiclass classifiers (Bury et al., 2021), our binary classifier also shows improved
377 performance (see Table 4). It requires significantly less computational power and time to train, as it only contains two
378 classes instead of several. For the length 500 classifier its performance metrics are better than the multiclass version
379 (see Table 4), owing to its more focused nature. The specific training only on fold bifurcation and null timeseries
380 prevents a ‘dilution’ of the classifier’s performance from the overlap of features of several different bifurcation types.
381 For the length 1500 classifier version the multiclass approach shows better performance indicating that longer
382 timeseries might benefit from classifiers that were trained on more than two classes. The disadvantage of the multiclass
383 models at shorter lengths potentially balances out at length 1500 due to learning from more representations, that allow
384 for better separation. However, as the geological record tends to be rather incomplete and a majority of geological
385 high-resolution timeseries are relatively short (less than 500 datapoints), our binary classifier constitutes a valuable
386 new tool for tipping point detection tailored to this kind of geological records.

387 Another advantage of our binary classifier over the existing multiclass classifiers is its ability to flag misclassified
388 transcritical bifurcations using an energy-based OOD detection approach. Transcritical and fold bifurcations are
389 mathematically very similar (Kuznetsov, 1998) and thus, current DL frameworks have issues picking up on their subtle
390 differences, often misclassifying transcritical bifurcations as fold (see our softmax confidence score results). This has
391 consequences for the interpretation of how abrupt and reversible an investigated transition is, an important
392 characteristic especially when studying and trying to mitigate Earth system tipping points. Transcritical bifurcations
393 are mostly gradual and easily reversible (Kuznetsov, 1998); misclassifying them as fold bifurcations leads to a high
394 rate of ‘false alarms’, in which timeseries are erroneously flagged as abrupt and irreversible tipping points. The energy-
395 based OOD detection approach implemented with our binary classifier is able to reduce the misclassification rate to
396 about 32% (length 500) and 34% (length 1500). This is a significant improvement over current multiclass classifiers,
397 that have a significantly higher rate of misclassification of about 47% (Bury et al., 2021). Thus, considering its
398 performance gain and lower rates of misclassification, our binary classifier together with the OOD approach constitute
399 an improvement on already existing DL multiclass classifiers.

400

401



402 4.2 Comparison of generic EWS and our binary classifier for the Cenozoic proxy records

403 We compare the results of our binary classifier with generic EWS from the literature for the different Cenozoic proxy
404 records, including the Paleocene/Eocene isotope record, the CENOGRID isotope record and the Eastern Mediterranean
405 sapropel record. Generic EWS from the literature pick up on several critical transitions in these records related to major
406 events associated with excursions in carbon and oxygen isotopes (Boettner et al., 2021; Setty et al., 2023) or the anoxia
407 proxy Mo (Hennekam et al., 2020a). We use our binary classifier to identify which of these transitions are fold
408 bifurcations and thus constitute an abrupt tipping point.

409 4.2.1 The Paleocene/Eocene isotope record

410 The PETM, ETM2 and ETM3 events are characterised by negative $\delta^{13}\text{C}$ as well as $\delta^{18}\text{O}$ excursions, that indicate strong
411 carbon cycle perturbations and warming (see Fig. 1), with the PETM being considered the strongest, followed by the
412 ETM2 and ETM3. For the PETM and ETM2 events generic EWS (autocorrelation and variance based) show a rising
413 trend (both autocorrelation and variance) in the $\delta^{13}\text{C}$ record, for the $\delta^{18}\text{O}$ isotope record only autocorrelation shows a
414 rising trend (Setty et al., 2023). A rise in variance and autocorrelation strongly implies an approaching transition
415 (associated with critical slowing down), a rise in autocorrelation only does not constitute a reliable bifurcation indicator
416 (Scheffer et al., 2009; Ditlevsen and Johnsen, 2010). The generic EWS results indicate approaching transitions for both
417 the PETM and ETM2 in the carbon isotope record, with an ambiguous signal for the oxygen isotope record. For the
418 PETM, our binary classifier length 1500 identifies fold bifurcations (tipping points) in both isotope records, with a
419 strong signal in $\delta^{13}\text{C}$ and to a lesser degree also in $\delta^{18}\text{O}$. However, our length 500 classifier does not recover any fold
420 bifurcations for the PETM. Depending on the choice of classifier, there is evidence for the PETM event being driven
421 by a tipping point in the carbon cycle (with the rise in temperature being a consequence rather than a driver of the
422 perturbation). Using our binary classifier, the ETM2 shows approaching fold bifurcations in both the $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$
423 isotope records, both with low to medium DL probabilities. This supports the ETM2 as being driven by perturbations
424 in the carbon cycle and temperature like the PETM, but with lower overall intensity. For the ETM3 generic EWS
425 results for both the $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ isotope records show no significant increase in autocorrelation and variance (except
426 for a rise in autocorrelation in the $\delta^{18}\text{O}$ record), implying no critical transitions associated with the ETM3 event (and
427 an ambiguous signal for the oxygen isotope record). Using our binary classifier we also do not recover any fold
428 bifurcation predictions; the weak ETM3 event is not associated with any tipping point. Our binary classifier generally
429 corroborates previous generic EWS results (for the PETM depending on the classifier length), that indicate approaching
430 transitions in the Paleocene/Eocene, and is able to identify these critical transitions as fold bifurcations (tipping points).

431 4.2.2 The CENOGRID isotope record

432 The Eocene/Oligocene Transition, the Oligocene/Miocene Transition and the Middle Miocene Transition are
433 associated with positive $\delta^{13}\text{C}$ as well as $\delta^{18}\text{O}$ excursions indicating a general cooling of the Earth system and ice volume
434 increase (see Fig. 2). For the Eocene/Oligocene Transition generic EWS do not indicate any combined rising trends in
435 autocorrelation and variance of the $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ record (only rising autocorrelation alone indicating an ambiguous
436 signal in these records) (Boettner et al., 2021). Our DL classifier returns null predictions for both isotope records,



437 corroborating the results of the generic EWS. For the Oligocene/Miocene Transition, generic EWS indicate critical
438 transitions in the $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ records based on rising trends in autocorrelation and variance for both records. Our
439 binary classifier length 500 also picks up on fold bifurcations in both isotope records with a higher DL probability for
440 the $\delta^{18}\text{O}$ record. However, using our length 1500 we do not recover any fold bifurcations. Based on the choice of
441 classifier, there is supportive evidence for the system going through a rapid and catastrophic tipping point at the
442 Oligocene/Miocene Transition driven mainly by perturbations in temperature (likely related to the development of a
443 stable Antarctic ice sheet). For the Middle Miocene Climate Transition, generic EWS do not return any combined rising
444 trends in autocorrelation and variance of the $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ record (except for a rise in the $\delta^{13}\text{C}$ autocorrelation),
445 indicating no approaching critical transitions (except for an ambiguous signal in the carbon isotope record). Our binary
446 classifier agrees with this by returning null predictions for both isotope records. Generally, our binary classifier
447 corroborates generic EWS results from the literature (for the Oligocene/Miocene transition depending on the choice of
448 classifier), that predict approaching transitions in the CENOGRID record, and is able to additionally identify these
449 transitions as fold bifurcations (tipping points).

450 **4.2.3 Comparison of DL classifier results and generic EWS results from the Eastern Mediterranean sapropel** 451 **record**

452 To assess the oxic-to-anoxic transitions associated with Mediterranean sapropels, the anoxia proxy Mo is used (see
453 Hennekam et al., 2020a). Generic EWS identify increases in variance of the Mo records prior to all anoxic events, all
454 within the 95% confidence level. Increases in autocorrelation are also found in all Mo records (with only the
455 autocorrelation results for S4 and S8 being within the 90% confidence level). Generic EWS indicate critical transitions
456 in all analysed sapropel events based on combined rising trends in variance and autocorrelation. Our binary classifier
457 supports this by returning fold predictions for all sapropels with consistently high DL predictions indicating that
458 sapropel events were actual tipping point events.

459 **4.3 Factors influencing classifier performance and results**

460 The performance of our DL classifiers on our empirical study systems can be influenced by a variety of factors.
461 Complex multivariate systems like the Paleocene/Eocene record and the CENOGRID record are more difficult to
462 analyze and can deliver conflicting or ambiguous DL results. They are influenced by various underlying coupled
463 processes, changing climatic boundary conditions and forcings (e.g. atmospheric CO_2 concentrations and existence of
464 polar ice sheets; see Westerhold et al., 2020). In addition, the CENOGRID record represents a composite record, based
465 on several cores from different latitudes and ocean basins. Every record from this composite could be influenced by
466 different forcings leading to varying tipping dynamics; combining these records might lead to ambiguous results in
467 our DL classifier. Generally, it seems that relatively simple systems with only few forcings like the Mediterranean
468 sapropel record (that is mostly driven by sea surface productivity and water column ventilation in a restricted basin;
469 see Rohling et al., 2015) tend to deliver more clear and unambiguous results supported by high DL probabilities. Our
470 analyses show, that results for timeseries of complex systems can be sensitive to the chosen classifier and the length
471 of the time interval studied prior to the bifurcation, which might lead to different results for the same timeseries
472 analysed with different classifier lengths. For example, for the PETM our binary classifier length 1500 returns fold



473 bifurcation predictions for both isotope records, whereas our binary classifier length 500 does not predict any
474 approaching transitions. In this case, the early warning and resilience loss of the system might have started long before
475 the bifurcation. Shortly before the bifurcation, other additional mechanisms than gradual resilience loss (which our DL
476 classifiers is trained to detect) could have been at play, ultimately pushing the system beyond its tipping point without
477 detectable early warning signals in the last 500 datapoints. Alternatively, the drivers of resilience loss related to the
478 PETM simply acted on longer timescales with lower frequencies and therefore need a longer classifier to be picked
479 up. We observe the opposite situation with the Oligocene/Miocene Transition, for which our length 1500 classifier
480 returns null predictions, whereas our length 500 classifier predicts fold bifurcations for both isotope records. Other
481 processes than gradual resilience loss (that our DL classifier is trained to detect) might act in the earlier parts of the
482 timeseries and obscure the signal leading to null predictions from our length 1500 DL classifier. Generally, our analyses
483 indicate, that there does not exist a ‘one size fits all’ classifier for all use cases, especially when considering complex
484 multivariate systems. A process-based understanding of the study system aids with the choice of classifier length. This
485 includes being aware of which forcings act on the investigated system and if they change over time. Through the right
486 choice of classifier length it is possible to include the forcings relevant for the tipping of the system. Additionally, it is
487 also important to have knowledge about the timeframes on which the relevant forcings act, in order to choose a suitable
488 classifier length capturing the relevant processes.

489 **4.4 Evaluation of OOD detection methods**

490 **4.4.1 OOD analyses on synthetically generated datasets**

491 Our softmax confidence score analyses reveals that 76% of the synthetically generated transcritical timeseries are
492 classified as fold bifurcations by our binary classifier. This indicates a high degree of overconfidence of our DL
493 classifier and a high rate of transcritical timeseries being misclassified as fold. Compared to a misclassification rate of
494 47% of already existing multiclass classifiers (Bury et al., 2021), this is comparably high and warrants alternative
495 mitigation strategies. Thus, we implemented our energy-based out-of-distribution (OOD) detection approach. This
496 approach is successful in separating in-distribution (fold) and out-of-distribution (transcritical) timeseries, especially
497 in the length 500 energy score diagram. Compared to using the softmax confidence scores, energy-based OOD
498 detection reduces the percentage of transcritical timeseries wrongly identified as fold bifurcations significantly from
499 76% to about 32% (length 500) and 34% (length 1500) respectively. This is also a significant improvement on the 47%
500 misclassification rate of already existing multiclass classifiers (Bury et al., 2021). We therefore propose the energy-
501 based OOD detection approach in combination with our binary classifier (especially the length 500 energy score
502 diagram and classifier) as an improved tool for the detection of fold bifurcations and flagging of potentially
503 misclassified transcritical transitions. The length 1500 energy score diagram does not show such a clear distinction
504 between ID and OOD as the length 500 energy score diagram (Fig. 5B). This might be caused by the smaller number
505 of synthetically generated timeseries it is based on; only 500 fold and transcritical timeseries each for length 1500
506 compared to 1250 timeseries each for length 500. This number was limited by the amount of the synthetically generated
507 timeseries in the original test dataset. This could potentially be improved upon in the future by generating more length
508 1500 test timeseries, in order to better separate ID and OOD in the length 1500 energy score diagram. However, as the



509 geological record is often incomplete and available timeseries rather short (below 500 datapoints), we focus on the
510 length 500 classifier and energy score diagram here.

511 We also identified noise as a potential cause for the bimodality of the ID and OOD energy score curves. There is a
512 significant difference in the noise levels, with timeseries with higher energy scores (above the energy threshold)
513 showing higher noise than the ones with lower energy scores (below the energy threshold). This means, that a higher
514 confidence in identification (indicated by higher energy scores) coincides with higher noise in the timeseries.
515 Generally, the detection of early warning signals is linked to the ‘noise’ portion of the timeseries. More stochastic
516 perturbation (noise) means that a system explores a larger portion of its local potential landscape (Scheffer et al., 2009;
517 Dakos et al., 2008) making it easier for a classifier to determine the type of bifurcation, thus being more confident and
518 returning higher energy scores. Low confidence energy scores are generally associated with less noisy data. In our
519 case, the binary classifier is less confident in identifying fold timeseries with little noise (below the energy threshold)
520 resulting in them being misclassified as transcritical. On the other hand, exceptionally noisy OOD transcritical
521 timeseries (above the energy threshold) are more similar to the ID fold bifurcations and have a higher chance of getting
522 misclassified as such. These findings indicate that our binary classifier is generally better at identifying fold
523 bifurcations, if they contain more noise. This can be an advantage over generic EWS approaches, that are typically not
524 well suited to handle noisy timeseries (Perretti and Munch, 2012; Dakos et al., 2015). Many empirical Earth systems
525 are open ‘noisy’ systems, governed by many (unknown) forcings and feedbacks (Dakos et al., 2015). Thus, for
526 detecting critical transitions in Earth systems with inherently higher amounts of noise, our new binary classifier
527 provides a powerful alternative to existing generic EWS.

528 **4.4.2 OOD analyses on empirical datasets**

529 Our energy-based OOD detection method to flag misclassified transcritical bifurcations identified all length 500
530 empirical timeseries previously classified as fold as genuine fold bifurcations. However, the length 1500 PETM
531 timeseries are flagged as misclassified transcritical bifurcations. As mentioned above, the length 1500 energy score
532 diagram is based on significantly less synthetically generated timeseries than the length 500 version and thus could
533 yield misleading results due to lacking data basis. Another reason could be, that the empirical PETM timeseries contain
534 only little amounts of noise and are therefore erroneously identified as OOD (transcritical) by our energy-based OOD
535 detection method. The empirical case of the PETM needs further investigation and testing to be able to unambiguously
536 identify whether this event was an actual catastrophic tipping point (fold bifurcation).

537

538 **5 Conclusions**

539 We developed and successfully tested a new deep learning classifier, focused on the detection of fold bifurcations in
540 the Earth system. Our binary classifier shows clear performance gains in comparison to generic EWS as well as existing
541 DL classifiers. We additionally implemented an energy-based OOD detection method that is able to significantly reduce
542 the amount of misclassified transcritical bifurcations. Comparing our DL approach with the results of previous studies
543 (using generic EWS) we are generally able to corroborate the existence of previously identified critical transitions. We



544 add to this by identifying most of the critical transitions from the empirical datasets as genuine fold bifurcation tipping
545 points using our binary classifier paired with energy-based OOD detection. There's scope for future investigation into
546 how exactly factors like classifier length or noise levels influence the results. Nevertheless, our binary classifier
547 together with the energy-based OOD detection method provides a powerful new tool for the reliable detection of
548 tipping points geared specifically towards Earth system applications.

549

550 **Code and data availability**

551 All data and code, as well as instructions for reproducing the results and figures in this manuscript are available on
552 Zenodo (Grohgan et al., 2026). To facilitate practical adoption, we will release an open-source implementation of our
553 classifier and out-of-distribution (OOD) detection method, along with detailed usage instructions and tutorials, in a
554 public GitHub repository upon publication.

555

556 **Author contribution**

557 M.G.: investigation, formal analysis, methodology, software, data curation, validation, visualization, writing – original
558 draft, writing – review and editing; T.M.B: methodology, software, data curation, writing – review and editing;
559 B.v.d.B.: formal analysis, writing – review and editing; G.J.R.: conceptualization, formal analysis, writing – review
560 and editing; R.H.: conceptualization, formal analysis, project administration, supervision, funding acquisition, writing
561 – review and editing.

562 All authors gave final approval for publication and agreed to be held accountable for the work performed therein.

563

564 **Competing interests**

565 The authors declare that they have no conflict of interests.

566

567 **Acknowledgements**

568 This work was supported by the Netherlands Organisation for Scientific Research (NWO) through the Vidi grant (grant
569 no. VI.Vidi.223.138) awarded to R.H. Additional support was provided by EMBRACER (Summit grant
570 SUMMIT.1.034), also financed by NWO.

571

572

573



574 References

- 575 Armstrong McKay, D. I., Staal, A., Abrams, J. F., Winkelmann, R., Sakschewski, B., Loriani, S., Fetzer, I., Cornell, S.
576 E., Rockström, J., and Lenton, T. M.: Exceeding 1.5°C global warming could trigger multiple climate tipping points,
577 *Science*, 377, eabn7950, <https://doi.org/10.1126/science.abn7950>, 2022.
- 578 Arnold, V. I.: Geometrical methods in the theory of ordinary differential equations, Springer, New York 2012.
- 579 Barnet, J. S. K., Littler, K., Westerhold, T., Kroon, D., Leng, M. J., Bailey, I., Röhl, U., and Zachos, J. C.: Late
580 Cretaceous-Early Paleogene stable isotope and coarse fraction record of ODP Site 208-1262, PANGAEA [dataset],
581 <https://doi.org/10.1594/PANGAEA.884588>, 2018.
- 582 Barnet, J. S. K., Littler, K., Westerhold, T., Kroon, D., Leng, M. J., Bailey, I., Röhl, U., and Zachos, J. C.: A High-
583 Fidelity Benthic Stable Isotope Record of Late Cretaceous–Early Eocene Climate Change and Carbon-Cycling,
584 *Paleoceanography and Paleoclimatology*, 34, 672-691, <https://doi.org/10.1029/2019PA003556>, 2019.
- 585 Boers, N.: Early-warning signals for Dansgaard-Oeschger events in a high-resolution ice core record, *Nature*
586 *Communications*, 9, 2556, <https://doi.org/10.1038/s41467-018-04881-7>, 2018.
- 587 Boers, N.: Observation-based early-warning signals for a collapse of the Atlantic Meridional Overturning Circulation,
588 *Nature Climate Change*, 11, 680-688, <https://doi.org/10.1038/s41558-021-01097-4>, 2021.
- 589 Boers, N. and Rypdal, M.: Critical slowing down suggests that the western Greenland Ice Sheet is close to a tipping
590 point, *Proceedings of the National Academy of Sciences*, 118, e2024192118,
591 <https://doi.org/10.1073/pnas.2024192118>, 2021.
- 592 Boettner, C., Klinghammer, G., Boers, N., Westerhold, T., and Marwan, N.: Early-warning signals for Cenozoic climate
593 transitions, *Quaternary Science Reviews*, 270, 107177, <https://doi.org/10.1016/j.quascirev.2021.107177>, 2021.
- 594 Bury, T. M.: ewstools: a Python package for early warning signals of bifurcations in time series data, *Journal of Open*
595 *Source Software*, 8, 5038, 2023.
- 596 Bury, T. M., Dylewsky, D., Bauch, C. T., Anand, M., Glass, L., Shrier, A., and Bub, G.: Predicting discrete-time
597 bifurcations with deep learning, *Nature Communications*, 14, 6331, <https://doi.org/10.1038/s41467-023-42020-z>,
598 2023.
- 599 Bury, T. M., Sujith, R. I., Pavithran, I., Scheffer, M., Lenton, T. M., Anand, M., and Bauch, C. T.: Deep learning for
600 early warning signals of tipping points, *Proceedings of the National Academy of Sciences*, 118, e2106140118,
601 <https://doi.org/10.1073/pnas.2106140118>, 2021.
- 602 Crawford, J. D.: Introduction to bifurcation theory, *Reviews of Modern Physics*, 63, 991-1037,
603 <https://doi.org/10.1103/RevModPhys.63.991>, 1991.
- 604 Dablander, F. and Bury, T. M.: Deep learning for tipping points: Preprocessing matters, *Proceedings of the National*
605 *Academy of Sciences*, 119, e2207720119, <https://doi.org/10.1073/pnas.2207720119>, 2022.
- 606 Dakos, V., Carpenter, S. R., van Nes, E. H., and Scheffer, M.: Resilience indicators: prospects and limitations for early
607 warnings of regime shifts, *Philosophical Transactions of the Royal Society B: Biological Sciences*, 370,
608 10.1098/rstb.2013.0263, 2015.
- 609 Dakos, V., Scheffer, M., van Nes, E. H., Brovkin, V., Petoukhov, V., and Held, H.: Slowing down as an early warning
610 signal for abrupt climate change, *Proceedings of the National Academy of Sciences*, 105, 14308-14312,
611 doi:10.1073/pnas.0802430105, 2008.
- 612 Deb, S., Sidheekh, S., Clements, C. F., Krishnan, N. C., and Dutta, P. S.: Machine learning methods trained on simple
613 models can predict critical transitions in complex natural systems, *Royal Society Open Science*, 9, 211475,
614 <https://doi.org/10.1098/rsos.211475>, 2022.
- 615 Ditlevsen, P. D. and Johnsen, S. J.: Tipping points: Early warning and wishful thinking, *Geophysical Research Letters*,
616 37, <https://doi.org/10.1029/2010GL044486>, 2010.
- 617 Drijfhout, S.: Competition between global warming and an abrupt collapse of the AMOC in Earth's energy imbalance,
618 *Scientific Reports*, 5, 14877, <https://doi.org/10.1038/srep14877>, 2015.
- 619 Galeotti, S., DeConto, R., Naish, T., Stocchi, P., Florindo, F., Pagani, M., Barrett, P., Bohaty, S. M., Lanci, L., Pollard,
620 D., Sandroni, S., Talarico, F. M., and Zachos, J. C.: Antarctic Ice Sheet variability across the Eocene-Oligocene
621 boundary climate transition, *Science*, 352, 76-80, <https://doi.org/10.1126/science.aab0669>, 2016.
- 622 Garbe, J., Albrecht, T., Levermann, A., Donges, J. F., and Winkelmann, R.: The hysteresis of the Antarctic Ice Sheet,
623 *Nature*, 585, 538-544, <https://doi.org/10.1038/s41586-020-2727-5>, 2020.
- 624 Grohgan, M., Bury, T. M., van der Bolt, B., Reichart, G.-J., and Hennekam, R.: An Earth system deep learning
625 classifier for tipping point detection: supplementary data and code, Zenodo [dataset],
626 <https://doi.org/10.5281/zenodo.19629897>, 2026.



- 627 Hawkins, E., Smith, R. S., Allison, L. C., Gregory, J. M., Woollings, T. J., Pohlmann, H., and de Cuevas, B.: Bistability
628 of the Atlantic overturning circulation in a global climate model and links to ocean freshwater transport, *Geophysical*
629 *Research Letters*, 38, <https://doi.org/10.1029/2011GL047208>, 2011.
- 630 Hendrycks, D. and Gimpel, K.: A baseline for detecting misclassified and out-of-distribution examples in neural
631 networks, arXiv preprint arXiv:1610.02136, 2016.
- 632 Hennekam, R., van der Bolt, B., van Nes, E. H., de Lange, G. J., Scheffer, M., and Reichart, G.-J.: Early-Warning
633 Signals for Marine Anoxic Events, *Geophysical Research Letters*, 47, e2020GL089183,
634 <https://doi.org/10.1029/2020GL089183>, 2020a.
- 635 Hennekam, R., Van der Bolt, B., Van Nes, E. H., de Lange, G. J., Scheffer, M., and Reichart, G.-J.: Calibrated XRF-
636 scanning data (mm resolution) and calibration data (ICP-OES and ICP-MS) for elements Al, Ba, Mo, Ti, and U in
637 Mediterranean cores MS21, MS66, and 64PE406E1, PANGAEA [dataset],
638 <https://doi.org/10.1594/PANGAEA.923197>, 2020b.
- 639 Jackson, L. C., Kahana, R., Graham, T., Ringer, M. A., Woollings, T., Mecking, J. V., and Wood, R. A.: Global and
640 European climate impacts of a slowdown of the AMOC in a high resolution GCM, *Climate Dynamics*, 45, 3299-3316,
641 <https://doi.org/10.1007/s00382-015-2540-2>, 2015.
- 642 Kéfi, S., Guttal, V., Brock, W. A., Carpenter, S. R., Ellison, A. M., Livina, V. N., Seekell, D. A., Scheffer, M., van Nes,
643 E. H., and Dakos, V.: Early Warning Signals of Ecological Transitions: Methods for Spatial Patterns, *PLOS ONE*, 9,
644 e92097, <https://doi.org/10.1371/journal.pone.0092097>, 2014.
- 645 Kermack, W. O. and McKendrick, A. G.: A contribution to the mathematical theory of epidemics, *Proceedings of the*
646 *Royal Society of London. Series A, Containing Papers of a Mathematical and Physical Character*, 115, 700-721,
647 <https://doi.org/10.1098/rspa.1927.0118>, 1927.
- 648 Kriegler, E., Hall, J. W., Held, H., Dawson, R., and Schellnhuber, H. J.: Imprecise probability assessment of tipping
649 points in the climate system, *Proceedings of the National Academy of Sciences*, 106, 5041-5046,
650 <https://doi.org/10.1073/pnas.0809117106>, 2009.
- 651 Kuznetsov, Y. A.: *Elements of applied bifurcation theory*, Springer 1998.
- 652 Lauretano, V., Zachos, J. C., and Lourens, L. J.: Benthic stable isotope record of ODP Site 208-1263, PANGAEA
653 [dataset], <https://doi.org/10.1594/PANGAEA.893894>, 2018a.
- 654 Lauretano, V., Zachos, J. C., and Lourens, L. J.: Orbitally Paced Carbon and Deep-Sea Temperature Changes at the
655 Peak of the Early Eocene Climatic Optimum, *Paleoceanography and Paleoclimatology*, 33, 1050-1065,
656 <https://doi.org/10.1029/2018PA003422>, 2018b.
- 657 Lear, C. H., Bailey, T. R., Pearson, P. N., Coxall, H. K., and Rosenthal, Y.: Cooling and ice growth across the Eocene-
658 Oligocene transition, *Geology*, 36, 251-254, <https://doi.org/10.1130/g24584a.1>, 2008.
- 659 Lenton, T. M., Held, H., Kriegler, E., Hall, J. W., Lucht, W., Rahmstorf, S., and Schellnhuber, H. J.: Tipping elements
660 in the Earth's climate system, *Proceedings of the National Academy of Sciences*, 105, 1786-1793,
661 <https://doi.org/10.1073/pnas.0705414105>, 2008.
- 662 Lenton, T. M., Milkoreit, M., Willcock, S., Abrams, J. F., Mc Kay, D. L. A., Buxton, J. E., Donges, J., Loriani, S.,
663 Wunderling, N., and Alkemade, F.: Global tipping points report 2025, 2025.
- 664 Liu, W., Wang, X., Owens, J., and Li, Y.: Energy-based out-of-distribution detection, *Advances in neural information*
665 *processing systems*, 33, 21464-21475, 2020.
- 666 Milkoreit, M., Hodbod, J., Baggio, J., Benessaiah, K., Calderón-Contreras, R., Donges, J. F., Mathias, J.-D., Rocha, J.
667 C., Schoon, M., and Werners, S. E.: Defining tipping points for social-ecological systems scholarship—an
668 interdisciplinary literature review, *Environmental Research Letters*, 13, 033005, <https://doi.org/10.1088/1748-9326/aaaa75>, 2018.
- 669 Perretti, C. T. and Munch, S. B.: Regime shift indicators fail under noise levels commonly observed in ecological
670 systems, *Ecological Applications*, 22, 1772-1779, <https://doi.org/10.1890/11-0161.1>, 2012.
- 672 Rahmstorf, S.: Bifurcations of the Atlantic thermohaline circulation in response to changes in the hydrological cycle,
673 *Nature*, 378, 145-149, <https://doi.org/10.1038/378145a0>, 1995.
- 674 Robinson, A., Calov, R., and Ganopolski, A.: Multistability and critical thresholds of the Greenland ice sheet, *Nature*
675 *Climate Change*, 2, 429-432, <https://doi.org/10.1038/nclimate1449>, 2012.
- 676 Rohling, E. J., Marino, G., and Grant, K. M.: Mediterranean climate and oceanography, and the periodic development
677 of anoxic events (sapropels), *Earth-Science Reviews*, 143, 62-97, <https://doi.org/10.1016/j.earscirev.2015.01.008>,
678 2015.
- 679 Scheffer, M., Bascompte, J., Brock, W. A., Brovkin, V., Carpenter, S. R., Dakos, V., Held, H., van Nes, E. H., Rietkerk,
680 M., and Sugihara, G.: Early-warning signals for critical transitions, *Nature*, 461, 53-59,
681 <https://doi.org/10.1038/nature08227>, 2009.



- 682 Setty, S., Cramwinckel, M. J., van Nes, E. H., van de Leemput, I. A., Dijkstra, H. A., Lourens, L. J., Scheffer, M., and
683 Sluijs, A.: Loss of Earth system resilience during early Eocene transient global warming events, *Science Advances*, 9,
684 eade5466, <https://doi.org/10.1126/sciadv.ade5466>, 2023.
- 685 Sexton, P. F., Norris, R. D., Wilson, P. A., Pälike, H., Westerhold, T., Röhl, U., Bolton, C. T., and Gibbs, S.: Eocene
686 global warming events driven by ventilation of oceanic dissolved organic carbon, *Nature*, 471, 349-352,
687 <https://doi.org/10.1038/nature09826>, 2011.
- 688 Seydel, R. U.: Practical bifurcation and stability analysis, Springer2009.
- 689 Shevenell, A. E., Kennett, J. P., and Lea, D. W.: Middle Miocene Southern Ocean Cooling and Antarctic Cryosphere
690 Expansion, *Science*, 305, 1766-1770, <https://doi.org/10.1126/science.1100061>, 2004.
- 691 Stommel, H.: Thermohaline Convection with Two Stable Regimes of Flow, *Tellus*, 13, 224-230,
692 <https://doi.org/10.1111/j.2153-3490.1961.tb00079.x>, 1961.
- 693 Van Breedam, J., Huybrechts, P., and Crucifix, M.: Hysteresis and orbital pacing of the early Cenozoic Antarctic ice
694 sheet, *Climate of the Past*, 19, 2551-2568, <https://doi.org/10.5194/cp-19-2551-2023>, 2023.
- 695 van Nes, E. H., Arani, B. M. S., Staal, A., van der Bolt, B., Flores, B. M., Bathiany, S., and Scheffer, M.: What Do You
696 Mean, ‘Tipping Point’?, *Trends in Ecology & Evolution*, 31, 902-904, <https://doi.org/10.1016/j.tree.2016.09.011>, 2016.
- 697 van Westen, R. M. and Dijkstra, H. A.: Asymmetry of AMOC Hysteresis in a State-Of-The-Art Global Climate Model,
698 *Geophysical Research Letters*, 50, e2023GL106088, <https://doi.org/10.1029/2023GL106088>, 2023.
- 699 van Westen, R. M., Kliphuis, M., and Dijkstra, H. A.: Physics-based early warning signal shows that AMOC is on
700 tipping course, *Science Advances*, 10, eadk1189, <https://doi.org/10.1126/sciadv.adk1189>, 2024.
- 701 van Westen, R. M., Vanderborcht, E., and Dijkstra, H. A.: A saddle-node bifurcation may be causing the AMOC
702 collapse in the Community Earth System Model, *Earth Syst. Dynam.*, 16, 2063-2085, [https://doi.org/10.5194/esd-16-](https://doi.org/10.5194/esd-16-2063-2025)
703 [2063-2025](https://doi.org/10.5194/esd-16-2063-2025), 2025.
- 704 Westerhold, T.: Cenozoic global reference benthic carbon and oxygen isotope dataset (CENOGRID) PANGAEA
705 [dataset], <https://doi.org/10.1594/PANGAEA.917503>, 2020.
- 706 Westerhold, T., Marwan, N., Drury, A. J., Liebrand, D., Agnini, C., Anagnostou, E., Barnet, J. S. K., Bohaty, S. M., De
707 Vleeschouwer, D., Florindo, F., Frederichs, T., Hodell, D. A., Holbourn, A. E., Kroon, D., Laurentano, V., Littler, K.,
708 Lourens, L. J., Lyle, M., Pälike, H., Röhl, U., Tian, J., Wilkens, R. H., Wilson, P. A., and Zachos, J. C.: An
709 astronomically dated record of Earth’s climate and its predictability over the last 66 million years, *Science*, 369, 1383-
710 1387, <https://doi.org/10.1126/science.aba6853>, 2020.
- 711 Wilson, G. S., Pekar, S. F., Naish, T. R., Passchier, S., and DeConto, R.: Chapter 9 The Oligocene–Miocene Boundary
712 – Antarctic Climate Response to Orbital Forcing, in: *Developments in Earth and Environmental Sciences*, edited by:
713 Florindo, F., and Siegert, M., Elsevier, 369-400, [https://doi.org/10.1016/S1571-9197\(08\)00009-8](https://doi.org/10.1016/S1571-9197(08)00009-8), 2008.
- 714 Wissel, C.: A universal law of the characteristic return time near thresholds, *Oecologia*, 65, 101-107,
715 <https://doi.org/10.1007/BF00384470>, 1984.
- 716 Wunderling, N., Donges, J. F., Kurths, J., and Winkelmann, R.: Interacting tipping elements increase risk of climate
717 domino effects under global warming, *Earth Syst. Dynam.*, 12, 601-619, <https://doi.org/10.5194/esd-12-601-2021>,
718 2021.
- 719 Zachos, J., Pagani, M., Sloan, L., Thomas, E., and Billups, K.: Trends, Rhythms, and Aberrations in Global Climate
720 65 Ma to Present, *Science*, 292, 686-693, <https://doi.org/10.1126/science.1059412>, 2001.
- 721 Zachos, J. C., Röhl, U., Schellenberg, S. A., Sluijs, A., Hodell, D. A., Kelly, D. C., Thomas, E., Nicolo, M., Raffi, I.,
722 Lourens, L. J., McCarren, H., and Kroon, D.: Rapid Acidification of the Ocean During the Paleocene-Eocene Thermal
723 Maximum, *Science*, 308, 1611-1615, <https://doi.org/10.1126/science.1109004>, 2005.

724