

1 **Distinguishing trends and shifts from memory in climate data**

2

3

4 Claudie Beaulieu^{1,2*} and Rebecca Killick³

5

6 ¹ Ocean Sciences Department, University of California, Santa Cruz, 95064, USA

7 ² Ocean and Earth Science, University of Southampton, Waterfront Campus, European
8 Way, Southampton, SO14 3ZH, UK

9 ³ Department of Mathematics and Statistics, Lancaster University, Lancaster, LA1 4YF, UK

10

11 * Corresponding author:

12 Claudie Beaulieu

13 Telephone: 831-459-5152

14 Email: beaulieu@ucsc.edu

15

16

17

18 Keywords: change-point detection, climate time-series, autocorrelation, memory, long-term
19 trend

20 **Abstract**

21 The detection of climate change and its attribution to the corresponding underlying processes
22 is challenging because signals such as trends and shifts are superposed on variability arising
23 from the memory within the climate system. Statistical methods used to characterize change
24 in time-series must be flexible enough to distinguish these components. Here we propose an
25 approach tailored to distinguish these different modes of change by fitting a series of models
26 and selecting the most suitable one according to an information criterion. The models involve
27 combinations of a constant mean or a trend superposed to a background of white-noise with
28 or without autocorrelation to characterize the memory, and is able to detect multiple change-
29 points in each model configuration. Through a simulation study on synthetic time-series the
30 approach is shown to be effective in distinguishing abrupt changes from trends and memory
31 by identifying the true number and timing of abrupt changes when they are present.
32 Furthermore, the proposed method is better performing than two commonly used approaches
33 for the detection of abrupt changes in climate time-series. Using this approach the so-called
34 “hiatus” in recent global mean surface warming fails to be detected as a shift in the rate of
35 temperature rise but is instead consistent with steady increase since the 1960s/1970s. Our
36 method also supports the hypothesis that the Pacific Decadal Oscillation behaves as a short-
37 memory process, rather than forced mean shifts as previously suggested. These examples
38 demonstrate the usefulness of the proposed approach for change detection and for avoiding
39 the most pervasive types of mistake in detection of climate change.

40 **1. Introduction**

41 The pace of climate change is not smooth; it varies year-to-year and decade-to-decade,
42 naturally. Climate records contain shifts or “abrupt changes” due to internal variability and
43 natural forcings (volcanic and solar) superimposed on the long-term anthropogenic climate
44 change trend (Fyfe et al. 2016; Lean and Rind 2009; Trenberth 2015). For example, the
45 global annual mean surface temperature (GMST) time-series exhibits periods of warming
46 separated by a long pause from approximately mid 1940s to mid 1970s (Kellogg 1993) and
47 potentially a second and shorter one, although highly debated, since the late 1990s/early
48 2000s (Drijfhout et al. 2014; Karl et al. 2015; Trenberth 2015; Trenberth and Fasullo 2013).
49 Whether this last so-called “hiatus” can be characterized as a slowdown in the rate of climate
50 change is the subject of active debate (Medhaug et al. 2017) and has led to a fast growing
51 number of scientific publications (Lewandowsky et al. 2016; Lewandowsky et al. 2015).
52 Discrepancies between the continued warming in models and apparent slowdown of warming
53 in observations since the late 1990s/early 2000s have been suggested to arise from
54 misrepresentations of forcing or natural variability in models (Huber and Knutti 2014; Meehl
55 et al. 2014; Risbey et al. 2014; Santer et al. 2014; Schmidt et al. 2014) or from data biases in
56 observations (Karl et al. 2015), and such change would unlikely be persistent (Knutson et al.
57 2016). However, few authors have addressed the problem from a statistical change detection
58 perspective (Cahill et al. 2015; Rahmstorf et al. 2017; Rajaratnam et al. 2015). From this
59 angle, the main question is whether the GMST trend has changed in the late 1990s/early
60 2000s and whether a significant slowdown of warming can be detected.

61 The Pacific Decadal Oscillation (PDO) has been suggested as a main driver of
62 variability in the GMST increase (Trenberth 2015), with its cold phases corresponding to
63 periods of paused warming and warm phases corresponding to GMST increase. The PDO has
64 also been suggested to be responsible for widespread ecosystem shifts in the North Pacific

65 with repercussions on the region's fisheries (Mantua et al. 1997) and drought effects of the El
66 Niño Southern oscillation (ENSO) (Wang et al. 2014). Whether PDO shifting patterns arise
67 from internal variability or from a forced bi-stable behavior has also triggered debate in the
68 literature over the last two decades (Mantua et al. 1997; Newman et al. 2016; Rodionov 2006;
69 Rudnick and Davis 2003), and has implications for its predictability.

70 Statistical approaches to characterize change in time-series behaving as a
71 superposition of several components such as long-term trends, shifts (i.e. either in the rate of
72 change or between two stable states) and internal variability, must be flexible enough to
73 distinguish these components. Internal variability is often characterized by a short-memory
74 process, in which the ocean and other slow components of the climate system (e.g. ice sheets)
75 respond slowly to random atmospheric forcing, producing climate variability at a longer time
76 scale than the white noise atmospheric weather. This mechanism is often referred to as “red
77 noise” in the climate literature (Frankignoul and Hasselmann 1977; Hasselmann 1976; Vallis
78 2010). Natural fluctuations caused by the internal memory can be large enough to mask the
79 long-term warming trend and create periods of apparent slowdown, possibly akin to a
80 “hiatus”, as well as exaggerate the warming trend for short periods, which implies risk for
81 ecosystems (Mustin et al. 2013). Long-term trends and shifts above that level of short-term
82 memory should represent natural or external forcings.

83 Climate science has typically put greater emphasis on statistical model interpretability
84 rather than flexibility because focus is more on a system-level understanding rather than
85 prediction of single events (Faghmous and Kumar 2014). Therefore, statistical approaches
86 used to quantify long-term change in climate time-series typically assume the change is linear
87 in time (Hartmann et al. 2013), and may not allow for all features described above in the
88 same model, thus leading to five possible misuses of statistics, which are illustrated in Fig. 1.

89 The first type of misuse can occur when characterizing GMST changes (Seidel and
90 Lanzante 2004), i.e. fitting a linear trend in presence of shifts in the mean or shifts in trend
91 (Fig. 1a), which can potentially bias the estimated rate of change. A series of alternative
92 piecewise linear models has been suggested to represent the GMST time-series including
93 periods of warming separated by a pause from the mid 1940s to 1970s (Seidel and Lanzante
94 2004). However, the performance of such piecewise models to characterize change in the
95 GMST depends on their ability to identify the timings separating the intervals of different
96 rates of warming. Advances in statistics allow identifying the timing of such changes in time-
97 series using change-point detection (Beaulieu et al. 2012; Reeves et al. 2007), and these
98 approaches have recently been used to analyze the GMST time-series by fitting piecewise
99 linear models to objectively detect the timing of changes in the rate of warming (Cahill et al.
100 2015; Rahmstorf et al. 2017; Ruggieri 2012). More commonly in climate studies, however,
101 change-point detection has been used to detect only shifts in the mean of a time-series, for
102 example by applying the STARS approach (Rodionov 2004). This often leads to the second
103 type of misuse (Fig. 1b): fitting shifts in the mean in presence of a background trend. Because
104 the null model of the STARS approach is a constant mean and not a secular trend, shifts in
105 the mean will tend to provide a better fit to the trend than a constant mean. As such, the
106 method typically interprets a trend as a “staircase” series of abrupt changes (Beaulieu et al.
107 2016). However, an approach based on model selection, allowing one to distinguish shifts in
108 the mean from a background trend, can prevent the problem of confusing different types of
109 signals as per the first and second misuses (Beaulieu et al. 2012; Reeves et al. 2007).

110 In addition to different types of signal that may be confused, internal variability may
111 also be misinterpreted as a forced signal, e.g. as a long-term trend or mean shifts (Fig. 1c-d).
112 Patterns created by the internal memory of the system are challenging signal detection in
113 climate time-series as they pose the risk to be misinterpreted as trends or shifts. The risk is

114 greater in presence of short records (Wunsch 1999). The short-term memory or “red noise” is
115 often represented by a first-order autocorrelation process, AR(1), and challenges signal
116 detection as the risk of false alarms is increased when using statistical techniques designed
117 for independent data (von Storch 1999; von Storch and Zwiers 1999). In trend detection, the
118 internal variability can be distinguished from a secular trend by fitting a regression model
119 containing a trend and AR(1) through generalized least squares (Chatfield 2003) or by
120 adjusting the sample size by the effective number of independent observations, which is
121 reduced in presence of autocorrelation (von Storch and Zwiers 1999), thus avoiding the third
122 misuse. As for detecting abrupt changes, some methods have proposed approaches to
123 distinguish change-points from autocorrelation using information criterion and Monte Carlo
124 methods (Beaulieu et al. 2012; Robbins et al. 2016), or pre-whitening of the time-series
125 (Robbins et al. 2016; Rodionov 2006; Serinaldi and Kilsby 2016; Wang 2008) to prevent
126 from the fourth misuse. Finally, as the natural variability is characterized by an AR(1)
127 process, it carries memory that offers short-term predictability. Forecasting a time-series
128 using a stationary AR(1) model when there is an underlying trend and/or shifts in the mean is
129 the fifth possible misuse (Fig. 1e) and will lead to poor predictions.

130 Our work is thus motivated by the need for distinguishing signals and internal
131 variability in climate and environmental time-series, which is fundamental to better
132 understanding their behavior and predicting future changes. We investigate the behavior of
133 the GMST and PDO time-series (Fig. 2) by developing an approach, which fits a series of
134 models to a time-series and identifies the most appropriate according to the Akaike
135 information criterion (AIC), which is twice the model likelihood penalized by the number of
136 parameters fitted. The models involve combinations of a constant mean or a trend, with a
137 background of white-noise or an AR(1) process, and include the possibility of change-points
138 in each model configuration so as to yield eight models in total (Fig. 3). When a model with

139 change-points is considered, the number is estimated using an optimal segmentation
140 algorithm (Killick et al. 2012). We refer to our approach as “Environmental time-series
141 change-point detection” (EnvCpt) and have also created software available as an R package
142 on the Comprehensive R Archive Network (CRAN) (Killick et al. 2016). Details on the
143 methodology are provided in the next section. We further demonstrate the appropriateness of
144 the methodology through a simulation experiment in which we apply EnvCpt to synthetic
145 time-series mimicking signals and noise observed in climate time-series such as the GMST
146 and the PDO. We compare our approach to two methodologies that have been used to
147 investigate change-points in the GMST and PDO time-series respectively. More specifically,
148 we compare EnvCpt with the STARS methodology (Rodionov 2004), which has been
149 designed to detect mean change-points and has been used to investigate change-points in the
150 PDO among many other applications in the climate and oceanography literature. We also
151 compare EnvCpt with a Bayesian linear regression multiple change-point detection method
152 (BMCpt), which has been used to investigate change-points in the GMST (Ruggieri 2012).

153

154 **2. Methods**

155 *a. Data*

156 We use five annual GMST datasets:

157 1) Met Office Hadley Centre and Climatic Research Unit surface temperature dataset
158 (HadCRUT4)

159 The HadCRUT4 dataset (version HadCRUT.4.5.0.0; available at
160 <http://www.metoffice.gov.uk/hadobs/hadcrut4/data/current/download.html>) (Morice et al.
161 2012) comprises sea surface temperatures (SST) from the Hadley Centre SST dataset version

162 3 (HadSST3; (Kennedy et al. 2011a, 2011b) and land surface temperatures from the Climatic
163 Research Unit version 4 (Jones et al. 2012). The dataset anomalies are relative to 1961-1990.

164 2) HadCRUT4 infilled by kriging (HadCRUT4krig)

165 We use a variation of the HadCRUT4 dataset, in which regions with no observations were
166 infilled by kriging, mainly across the Arctic, Antarctic, parts of Africa and other small areas
167 (Cowtan and Way 2014); available at [http://www-](http://www-users.york.ac.uk/~kdc3/papers/coverage2013/series.html)
168 [users.york.ac.uk/~kdc3/papers/coverage2013/series.html](http://www-users.york.ac.uk/~kdc3/papers/coverage2013/series.html)). The reference period for the
169 anomalies is the same as for HadCRUT4.

170 3) Merged Land–Ocean Surface Temperature Analysis (MLOST)

171 The MLOST dataset from the National Oceanic and Atmospheric Administration National
172 Centers for Environmental Information (Smith et al. 2008; Vose et al. 2012; available at
173 <https://www.ncdc.noaa.gov/cag/time-series/global>) combines land air temperatures from the
174 Global Historical Climatology Network version 3.3.0 (GHCNv3.3.0) and the Extended
175 Reconstructed Sea Surface Temperature version 4 (ERSST.v4) (Huang et al. 2015; Liu et al.
176 2015). The anomalies are with respect to the 1971-2000 period.

177 4) Goddard Institute for Space Studies Surface Temperature Analysis (GISTEMP)

178 The GISTEMP dataset also combines land and SST temperatures from GHCNv3.3.0 and
179 ERSSTv4, but also includes the Scientific Committee on Antarctic Research (SCAR) stations
180 over Antarctica (Hansen et al. 2010) available at <http://data.giss.nasa.gov/gistemp/>). The
181 anomalies are relative to 1951-1980.

182 5) Berkeley Earth Surface Temperatures (BEST)

183 The BEST dataset (Rohde et al. 2013; available at <http://berkeleyearth.org/data/>) uses SST
184 derived from HadSST3 combined with air temperatures from CRUTEM4 and stations from
185 the GHCN network. Anomalies are given with respect to 1961-1990.

186 We use the HadCRUT4, HadCRUT4krig and BEST annual GMST datasets from 1850-2016
187 and the MLOST and GISTEMP annual GMST datasets from 1880-2016 (Figure 2). These
188 datasets share core common observations, but have been processed, bias-corrected and
189 interpolated independently (Jones and Kennedy 2017; Jones 2016).

190 The PDO dataset used was derived as the leading principal component of monthly sea surface
191 temperature in the North Pacific (downloaded from:
192 <http://jisao.washington.edu/pdo/PDO.latest>) (Mantua et al. 1997; Zhang et al. 1997). Annual
193 means from 1900-2016 were calculated from the monthly values as a mean from January to
194 December for each year, and presented in Figure 2.

195 *b. EnvCpt description*

196 EnvCpt fits eight models often used to represent climate and environmental time-
197 series and selects which one provides the best fit to represent the time series. The simplest
198 models for the time-series assume that the series is well represented by either a constant mean
199 or a linear trend in addition to a background white noise. These simple models are also fitted
200 superposed to an AR(1), leading to four types of models without change-points. Then,
201 models including change-points in all model parameters (mean or trend, variance and
202 autocorrelation) are also fitted, leading to a total of eight models that are described below.

203 1) a constant mean (Mean)

$$204 \quad y_t = \mu + e_t \quad (1)$$

205 where y_t represents the time-series, t is the time, μ is the mean and e_t is the white-noise

206 errors, which are independent and identically distributed following a Normal with a mean of
 207 zero and variance σ^2 .

208 2) a constant mean with first-order autocorrelation (Mean + AR(1))

$$209 \quad y_t = \mu + \varphi y_{t-1} + e_t \quad (2)$$

210 where φ is the first-order autocorrelation coefficient.

211 3) a linear trend (Trend)

$$212 \quad y_t = \lambda + \beta t + e_t \quad (3)$$

213 where λ and β represent the intercept and trend parameters, respectively.

214 4) a linear trend with first-order autocorrelation (Trend + AR(1))

$$215 \quad y_t = \lambda + \beta t + \varphi y_{t-1} + e_t \quad (4)$$

216 5) multiple change-points in the mean

$$217 \quad y_t = \begin{cases} \mu_1 + e_t & t \leq c_1 \\ \mu_2 + e_t & c_1 < t \leq c_2 \\ \vdots & \\ \mu_m + e_t & c_{m-1} < t \leq n \end{cases} \quad (5)$$

218 where μ_1, \dots, μ_m represent the mean of each of the m -segments with variance $\sigma_1^2, \dots, \sigma_m^2$
 219 respectively, c_1, \dots, c_{m-1} the timing of the change-points between segments and n is the length
 220 of the time-series.

221 6) multiple change-points in the mean and first-order autocorrelation

$$222 \quad y_t = \begin{cases} \mu_1 + \varphi_1 y_{t-1} + e_t & t \leq c_1 \\ \mu_2 + \varphi_2 y_{t-1} + e_t & c_1 < t \leq c_2 \\ \vdots & \\ \mu_m + \varphi_m y_{t-1} + e_t & c_{m-1} < t \leq n \end{cases} \quad (6)$$

223 where $\varphi_1, \dots, \varphi_m$ represent the autocorrelation in each segment.

224 7) a trend with multiple change-points in the regression parameters

$$225 \quad y_t = \begin{cases} \lambda_1 + \beta_1 t + e_t & t \leq c_1 \\ \lambda_2 + \beta_2 t + e_t & c_1 < t \leq c_2 \\ \vdots & \\ \lambda_m + \beta_m t + e_t & c_{m-1} < t \leq n \end{cases} \quad (7)$$

226 where $\lambda_1, \dots, \lambda_m$ and β_1, \dots, β_m represent the intercept and trend in each segment.

227 8) a trend with multiple change-points in the regression parameters and first-order
228 autocorrelation (Trend cpt + AR(1))

$$229 \quad y_t = \begin{cases} \lambda_1 + \beta_1 t + \varphi_1 y_{t-1} + e_t & t \leq c_1 \\ \lambda_2 + \beta_2 t + \varphi_2 y_{t-1} + e_t & c_1 < t \leq c_2 \\ \vdots & \\ \lambda_m + \beta_m t + \varphi_m y_{t-1} + e_t & c_{m-1} < t \leq n \end{cases} \quad (8)$$

230 The theoretical parameter ranges are real numbers for the means, trends and intercepts,
231 positive real numbers for the variances, $[-1,1]$ for first-order autocorrelation coefficients and
232 $[p, n-p]$ for the change-point timings with p parameters in the model form. The methodology
233 considers all possible parameters and number of changes across the 8 models.

234 Each model is fitted according to maximum likelihood estimation. For the change-
235 point models, we find the number and location of change-points using the Pruned Exact
236 Linear Time (PELT) algorithm (Killick et al. 2012), which identifies change-points by
237 performing an exact search considering all options for any possible number of changes
238 (varying from 1 to the maximum number of change-points given the set minimum segment
239 length). The search strategy is exact with a computational cost that is linear in the number of
240 data points. The PELT method is used in combination with the modified Bayesian
241 information criterion (MBIC) as the penalty function (Zhang and Siegmund 2007) to select
242 the optimal number of change-points, as this approach balances the overall fit against the
243 length of each segment. Hence it naturally guards against small segments unless it produces a

244 significantly improved fit. The PELT methodology may choose no change-point as the best
 245 model in which it reduces to the same likelihood as the no change equivalent model. The
 246 model selection is automated using the Akaike information criterion (AIC), which penalizes
 247 the model likelihood by the number of parameters fitted for each model considered (Akaike
 248 1974). The EnvCpt package provides the likelihood and number of parameters fitted for each
 249 model. As such, any other criteria or metric based on the likelihood can be used for the model
 250 selection. However, we use the MBIC for determining change-points as the AIC has been
 251 shown to systematically overestimate the number of changes (Haynes et al. 2017). The
 252 pseudo algorithm for EnvCpt and additional details about PELT are presented in Appendix A.

253 The best model is selected as the one with the smallest AIC. While the choice
 254 according to the minimum AIC does not provide a measure of uncertainty, the AIC
 255 differences (Δ_i) between the best model and the remaining models can be used to evaluate
 256 plausibility of the models fitted:

$$257 \quad \Delta_i = AIC_i - AIC_{min} \quad (9)$$

258 where i denotes the models fitted ($i=1, \dots, 8$). The larger the difference, the less plausible a
 259 model is, given the data and models considered (Burnham and Anderson 2002). As a rule of
 260 thumb, a Δ_i of 0-2 provides substantial support for model i , while Δ_i of 4-7 has considerably
 261 less support, and essentially none if the difference is larger than 10 (Burnham and Anderson
 262 2002). While comparing the differences to a rule of thumb is useful to identify a subset of
 263 models at play, we can also quantify the plausibility of the models fitted given the data using
 264 Akaike weights:

$$265 \quad w_i = \frac{\exp(-0.5 \cdot \Delta_i)}{\sum_{r=1}^8 \exp(-0.5 \cdot \Delta_r)} \quad (10)$$

266 The weights, w_i , represent the evidence in favor of model i being the best model given the
267 data and the set of eight models fitted.

268 *c. Simulation of synthetic series*

269 Synthetic series mimicking typical features observed in GMST and PDO time series
270 issued from the eight general models described in the previous section were generated to
271 assess the performance of EnvCpt. We generated a set of synthetic series inspired by the
272 GMST record with a total of 166 years that corresponds to the four models including a trend
273 component fitted to the GMST (Fig. 3a) with a) a long-term trend, b) a long-term trend with
274 first-order autocorrelation, c) a trend with three change-points in 1906, 1945 and 1976, and d)
275 one change-point in the trend and autocorrelation in 1962. We also generated synthetic time-
276 series inspired by the PDO with a length of 116 years to represent the competing models
277 suggested to characterize the PDO behavior: a) mean change-points in 1948 and 1976 with or
278 without a background of AR(1) (Rodionov 2004, 2006) and b) first-order autocorrelation
279 model (Newman et al. 2016). For completeness, the constant mean model used here
280 represents a “null” model for the two hypotheses. Figure 4 presents the eight cases of
281 synthetic series generated to mimic the GMST and PDO. The specific parameters used to
282 simulate the synthetic series are presented in Appendix A (Table A1). For each category, a
283 total number of 1,000 synthetic series were generated and analyzed.

284 *d. Comparison with STARS*

285 We compare our approach to STARS (Rodionov 2004, 2006) using the code available
286 from <http://www.climatelogic.com/download>. This approach has been used previously to
287 investigate the presence of mean shifts in the PDO (Rodionov 2004, 2006). STARS uses a
288 binary segmentation algorithm that identifies changes sequentially. As such, this procedure
289 finds the most likely change-point, then splits the data at the change if it is significant, and

290 searches for further changes in each segment. This procedure is repeated iteratively until no
291 more changes are detected or the segments are becoming smaller than the set minimum
292 segment length. The decision rule for the presence of change-points is based on a t-test
293 between segments (Rodionov 2004). A minimum segment length default of 10 observations
294 and a critical level of 5% were used in the present study. Thus we set the same default
295 minimum segment length with EnvCpt to carry out the simulations, although other options
296 can be used. The STARS methodology is developed to detect shifts in the mean, however we
297 present results for all considered models to demonstrate the errors produced when trends are
298 not accounted for within the model. Furthermore, STARS is not originally designed to handle
299 autocorrelation, and pre-whitening of the time-series has been suggested when its presence is
300 suspected (Rodionov 2006). Thus, we also applied STARS with two pre-whitening
301 approaches after some parameter tuning (Appendix C). The results obtained after pre-
302 whitening are presented in Appendix D.

303 *e. Comparison with BMCpt*

304 We also compare our approach to a Bayesian identification of multiple change-points
305 in a regression model (BMCpt), which has been used to investigate the presence of change-
306 points in the GMST (Ruggieri 2012). We use the code made freely available from
307 <http://mathcs.holycross.edu/~eruggier/software.html>. This approach allows for the detection
308 of changes in the parameters of a regression model and thus can detect changes in the mean,
309 trend and/or variance. The exact solution to the multiple change-point detection is obtained
310 using dynamic programming recursions. Here we use a minimum segment length between
311 two shifts of 10, the same as used for EnvCpt and STARS. This approach necessitates setting
312 several other parameters, which are chosen as per the recommendations in Ruggieri (2012)
313 and are described in Appendix B. The hyper-parameters for the variance prior are optimized,
314 as these have an effect on the number of change-points detected (Fig. A1; Appendix B).

315 BMCpt is also designed to fit a regression model with independent residuals. Thus, we also
316 apply it to the models with AR(1) after pre-whitening. Again, the choice of pre-whitening
317 parameters is determined by optimizing them to give the best performance and is presented in
318 Appendix C.

319

320 **3. Results**

321 *a. Analysis of the GMST and PDO time-series*

322 The eight EnvCpt models are fitted to the GMST datasets and the PDO in Fig. 3.
323 Table 1 presents the AIC differences for each model and their respective weights. For most
324 datasets, the evidence for the Trend cpt + AR(1) model is strong, with probabilities of 1 for
325 BEST, MLOST and GISTEMP, respectively (Table 1). For these three datasets, none of the
326 seven other models are considered plausible ($\Delta_i > 10; w_i = 0; i = 1, \dots, 7$). The
327 HadCRUT4krig dataset reveals more uncertainty, with substantial evidence for both the
328 Trend cpt + AR(1) and the Trend cpt models ($\Delta_i < 2; i = 7, 8$), but a higher probability for
329 the Trend cpt + AR(1) model (0.68 for Trend cpt + AR(1) as opposed to 0.32 for Trend cpt;
330 Table 1). On the opposite, for the HadCRUT4 dataset the best model is Trend cpt with a
331 probability of 0.98, while there is limited evidence for the Trend cpt + AR(1) model
332 (probability of 0.02).

333 For most GMST datasets, the best model fit has one change-point in both the trend
334 and autocorrelation (Trend cpt + AR(1)) in 1962 or 1972 depending on the source of the
335 GMST data (Fig. 3b-e; Table 1). At that time, the rate of warming increases and is
336 accompanied by a whitening of the GMST, i.e. the AR(1) weakens. The trend and AR(1)
337 parameters associated with this fit are presented in Table 2. The competing model (Trend cpt)
338 exhibits a flat mean until 1906, which was followed by a warming period until 1945, then

339 another period of minimal temperature change that lasted until 1977, followed by a warming
340 trend until now (Fig. 3a-b). It must be noted that all models fitted are valid if their underlying
341 assumptions of normality and independence of the residuals are met. Overall, these
342 assumptions are verified under the Trend cpt + AR(1) fit, but not under the Trend cpt model
343 (Figs A5-A6, Table A2; Appendix E). This further validates a background AR(1) and the
344 occurrence of one change-point in the GMST in 1962 or 1972, as opposed to several changes.
345 The GMST has also been suggested to follow an AR(2) model previously (Karl et al. 2000).
346 We find that while two datasets indicate a potential AR(2) structure in the residuals (Fig.
347 A6a-b; Appendix E), the fits are valid with an AR(1) (Fig. A5, Table A2; Appendix E).
348 Furthermore, an AR(2) does not seem to improve the likelihood of the model enough to be
349 worth including as all models with an AR(2) lead to substantially higher AIC (Table A2;
350 Appendix E).

351 The only model detecting a change-point in the late 1990s/early 2000s is the
352 “staircase” model (Mean cpt), for which there is essentially no evidence ($w_5 = 0$), given the
353 datasets and other models considered (Fig. 3a-e). As such, this result suggests that the most
354 recent “hiatus” does not emerge as a global signal, but rather indicates that the GMST rate of
355 change has remained approximately constant (linear) since the 1960s/1970s with some
356 fluctuations arising from the memory in the system.

357 As for the PDO, the best fitting model is a constant mean and autocorrelation (Mean +
358 AR(1)) with a probability of 0.56 (Table 1; Fig. 3f), and has valid underlying assumptions
359 (Fig. A7; Table A2). None of the models including change-points are considered at play, as
360 either no change-points are detected (Mean cpt + AR(1) and Trend cpt + AR(1)) or they are
361 associated with large AIC differences (Table 1). The Trend + AR(1) model is the only
362 competing model ($\Delta_4 = 1.1$; $w_4 = 0.44$), unveiling some uncertainty about the best way to

363 characterize PDO behavior. However, models including a trend would be counterintuitive to
364 represent PDO behavior (Newman et al. 2016).

365 *b. Simulation study*

366 EnvCpt was also applied to the eight different sets of synthetic series generated. To
367 emphasize the flexibility of the methodology developed, we compare it with two other
368 approaches both detailed in Methods. It must be noted that EnvCpt is developed to
369 distinguish all combinations of trends, change-points and autocorrelation, and thus we expect
370 it to overall outperform BMCpt and STARS, which are both designed for more specific
371 features. Specifically, BMCpt was developed to detect changes in a linear regression model,
372 and it should thus perform similarly to EnvCpt in presence of a constant mean or trend, with
373 or without change-points (cases Mean, Mean cpt, Trend and Trend cpt). Correspondingly,
374 STARS was developed to detect mean shifts only and should be performing in the simulation
375 scenario cases Mean and Mean cpt. Neither STARS nor BMCpt were originally designed to
376 handle a background of autocorrelation. To work around that limitation we also apply the
377 methods on the synthetic series with AR(1) after pre-whitening, which necessitates some
378 parameter tuning (see Appendix D).

379 Fig. 5 presents the number of shifts detected by EnvCpt, STARS and BMCpt in each
380 simulation case. The results demonstrate that EnvCpt correctly identifies the number of
381 change-points at a higher frequency than STARS and BMCpt in most synthetic series,
382 although BMCpt is equivalent in half of the cases. In presence of a trend only, both EnvCpt
383 and BMCpt succeed at identifying no change (Fig. 5a). However, in presence of three trend
384 change-points (Fig. 5c) EnvCpt detects the three shifts at the highest frequency while BMCpt
385 tends to interpret them as two shifts instead. The rate of false detection with BMCpt increases
386 in presence of autocorrelation (Fig. 5b), illustrating misuse 3. In the simulation case Trend

387 $\text{cpt} + \text{AR}(1)$, EnvCpt and BMCpt are equivalent (Fig. 5d) even though BMCpt is not
388 designed to handle autocorrelation. We attribute this result to the fact that BMCpt can detect
389 changes in the variance, thus interpreting the changing $\text{AR}(1)$ here as a change in variance.
390 Finally, in presence of mean shifts (cases Mean cpt and $\text{Mean cpt} + \text{AR}(1)$), BMCpt tends to
391 detect fewer shifts than the true number of change-points (Fig. 5g-h). Indeed, when using a
392 change-point approach fitting a piecewise linear regression model in presence of mean shifts
393 only, consecutive “staircase” mean shifts may be interpreted as a trend as per misuse 1. Pre-
394 whitening reduces the rate of false detection by BMCpt in the $\text{Trend} + \text{AR}(1)$ scenario, but
395 also diminishes the power of detection for the $\text{Trend cpt} + \text{AR}(1)$ and $\text{Mean cpt} + \text{AR}(1)$
396 cases (Fig. A3; Appendix D).

397 STARS tends to overestimate the number of change-points and frequently
398 misidentifies an underlying trend as a series of shifts, illustrating misuse 2 (Fig. 5a-d). In the
399 cases of a constant mean or change-points in the mean, STARS should be equivalent to
400 EnvCpt , but tends to detect additional spurious shifts (Fig. 5e,g). This is particularly
401 surprising for the Mean case (Fig. 5e), as the STARS methodology should be able to return a
402 no change model in this case, but rather detects changes in over 34% of the series. However,
403 although a 5% critical level is used when multiple shifts are present this does not correspond
404 to a 5% critical level for the overall segmentation given that the test is applied repetitively.
405 Approaches based on a maximal type t-test or F-test, which accounts for the fact that the test
406 statistic is calculated for each potential change-point timing in the time-series, reduce false
407 alarms to the expected level (Lund and Reeves 2002; Wang et al. 2007). The tendency for
408 spurious detection with STARS is aggravated in presence of autocorrelation (Fig. 5f), where
409 STARS detects changes in 96% of the series when none should be detected, illustrating
410 misuse 4. The rate of false detection is reduced with pre-whitening and the detection power
411 improved for the $\text{Mean} + \text{AR}(1)$ and $\text{Mean cpt} + \text{AR}(1)$ cases (Fig. A3; Appendix D).

412 Whilst the number of positive and false-positive changes detected by a given model
413 provides a picture of the performance, it does not indicate whether the change-points are
414 correctly localized in the time-series. Fig. 6 presents density estimates of the locations of the
415 identified change-points for synthetic series that were generated with change-points. This
416 again demonstrates that EnvCpt outperforms STARS and BMCpt overall. EnvCpt clearly
417 identifies the location of the trend change-points, while both BMCpt and STARS tend to
418 detect spurious changes between the true change-points (Fig. 6a), especially towards the end
419 of the series with STARS (Fig. 6a-b,d). The three methods are equivalent in detecting the
420 location of the mean change-points (Fig. 6c). It must be noted that the height of the density
421 peaks may suggest that BMCpt is better performing in the Mean cpt + AR(1) scenario, but
422 this is due to fewer changes being detected with this approach (Fig. 5h). The density and
423 number of change-points should be considered together.

424

425 **Discussion**

426 Our results suggest that the GMST rate of change has changed once in 1962 or 1972
427 and has remained approximately constant since then with fluctuations due to the presence of
428 memory in the system. Furthermore, we find that the GMST is “whitening” around that time,
429 i.e. the AR(1) parameter weakens. This result is consistent across most datasets with high
430 evidence (Table 1). Our GMST characterization is different from previous parametric
431 change-point analysis of the global temperature record (Cahill et al. 2015; Rahmstorf et al.
432 2017; Ruggieri 2012) that suggested the presence of three change-points in the GMST rate of
433 warming in the 1900s, 1940s and 1970s. The main difference lies in the treatment of
434 autocorrelation: our approach formally takes into account the autocorrelation by the means of
435 an AR(1). Indeed, the optimal fit of the Trend cpt model for the HADCRUT4 dataset (Fig.

436 3a), which does not take account of AR(1), detects three change-points as in previous studies.
437 However, autocorrelation is present in the residuals such that the underlying assumption of
438 independent residuals is violated under the Trend cpt model. The timings of change-points
439 under this model setting (Trend cpt) are not consistent across all GMST datasets, signaling
440 additional uncertainty. If the BIC is used to select the best model instead of the AIC, the
441 Trend cpt + AR(1) model is selected for all datasets (Table A4). We therefore argue that the
442 Trend cpt model should not be used without AR(1) to characterize the GMST. The GMST
443 has also been suggested to follow an AR(2) model previously (Karl et al. 2000). Here we find
444 that an AR(2) does not improve the likelihood of the model enough to be worth including as
445 the noise term (Table A2; Appendix E). Previous work has also suggested the presence of
446 long-term memory in surface temperature records (e.g. Franzke 2012; Løvsletten and Rypdal
447 2016), as opposed to the short-term memory detected here. In presence of long-term memory,
448 the autocorrelation function will not decay exponentially as observed here, but rather decays
449 as a power law such that it does not reach zero (Yuan et al. 2015). While we do not find long-
450 term memory in the residuals of the five GMST records analyzed here, we acknowledge that
451 its potential presence presents a risk to misinterpret it as a trend or an abrupt change with
452 EnvCpt, but longer records will be needed to make this distinction (Poppick et al. 2017).

453 Consequently, our results suggest that the change-points previously detected in the
454 1900s and 1940s may not be unusual given the background memory. These timings also
455 coincide with the period of highest uncertainty in SST measurements due to corrections
456 applied to account for changes of instrumentation (Jones 2016; Kent et al. 2017; Thompson et
457 al. 2008). Despite different results due to different change-point detection approach, we do
458 agree with previous studies (Cahill et al. 2015; Rahmstorf et al. 2017; Ruggieri 2012) that the
459 most recent “hiatus” in GMST does not emerge as a global signal, regardless of whether or
460 not AR(1) is considered. Hence, the only model fitted that contains a change-point in the late

461 1990s/early 2000s is a “staircase” in the GMST (Mean cpt) and that model fit is rendered
462 unlikely by its large AIC values (Fig. 3).

463 It must be noted that the five datasets employed in this study are not independent:
464 they all use in part the same input data for the land and ocean but employ different
465 methodologies for correcting biases and inhomogeneities and for interpolating (Jones 2016).
466 As such, the similar results obtained with the five datasets do not provide independent pieces
467 of evidence that a change-point took place in 1962 or 1972, but rather provides a measure of
468 the uncertainty arising from the different approaches used to create these datasets.

469 To our knowledge, the whitening of the GMST has not been described in previous
470 studies because methodologies able to detect shifts in the autocorrelation, such as EnvCpt,
471 have not been applied to GMST datasets before. The sudden decrease in memory detected
472 here could be due to changes in SST measurements, as the timing marks the start of a period
473 of SST measurements obtained from a more diverse observing fleet and reduced bias (Kent et
474 al. 2017; Thompson et al. 2008). Future studies should investigate the regions responsible for
475 the change-point in GMST and investigate the underlying causes.

476 As for the PDO, we show that a model with a flat mean and first-order autocorrelation
477 provides the best fit (Fig. 3f), which is in agreement with previous studies (Newman et al.
478 2016; Rudnick and Davis 2003). Conversely, a previous study has interpreted the PDO as a
479 series of shifts in the mean in the 1940s and 1970s, superposed to an AR(1) (Rodionov 2006),
480 which was taken as support for the hypothesis of a bi-stable behavior. When focusing on a
481 shorter period of time, the 1970s shift was also suggested to emerge from the background of
482 autocorrelation, although the authors questioned the robustness of this result and emphasized
483 the need of a methodology such as the one presented here (Beaulieu et al. 2016). Our new
484 methodology formally compares the two statistical representations (AR(1) process vs bi-

485 stability with mean shifts) of the PDO by considering them objectively, and we conclude that
486 it is best modeled as autocorrelation only, without shifts. This result is consistent if the BIC is
487 used to select the best model instead of the AIC (Table A4). Memory in the PDO can offer
488 short-term predictability a few years ahead, depending on the strength of the first-order
489 autocorrelation. Specifically, the first-order autocorrelation of 0.55 in the PDO time-series
490 analyzed here translates into a decorrelation time of 3.5 years (von Storch and Zwiers 1999)
491 after which the current PDO value will be “forgotten”. This predictability could be key for
492 management, as PDO patterns have widespread repercussions and have been suggested to be
493 responsible for ecosystem regime shifts in the North Pacific and regional droughts (Mantua et
494 al. 1997; Wang et al. 2014). More recently, it has been suggested that the PDO is “reddening”
495 at the monthly timescale, i.e. the AR(1) is increasing as a sign of critical slowing down
496 (Boulton and Lenton 2015; Lenton et al. 2017). We do not detect this feature here, but this is
497 not surprising since our approach is not designed to detect a trend in autocorrelation and has
498 been applied at the annual timescale.

499 As the PDO and GMST records become longer, the best fitting model may change.
500 More precisely, EnvCpt is expected to select the true underlying model and detect changes
501 more accurately as the number of observations increase (Killick et al. 2012).

502 The simulation study demonstrates the advantage of a single comprehensive method
503 to avoid five misuses of statistics in analyzing climate time-series. Our approach reduces the
504 number of pre-assumptions about the presence of trends, shifts and autocorrelation in the
505 time-series. In eight cases of synthetic series mimicking features observed in the GMST and
506 the PDO, our approach shows high skill in selecting the correct number of change-points in
507 mean and slope, and to locate the change-points correctly when present. A drawback is that
508 our conclusions are limited to the synthetic series generated for our simulation study.
509 However, previous simulation studies of change-point detection techniques on synthetic

510 series with shifts having a random timing and magnitude have been carried out before, and
511 revealed expected features that are common to most techniques. First, the signal-to-noise
512 ratio matters the most, i.e. a shift with a large magnitude compared to the background noise
513 has a higher hit rate (Beaulieu et al. 2012; Beaulieu et al. 2008; Reeves et al. 2007; Wang et
514 al. 2010). Second, false alarms occur more often at the beginning or end of the time-series
515 (Beaulieu et al. 2012). Third, successive shifts that are near in time tend to be more difficult
516 to detect, especially if the magnitudes have the same sign (e.g. an increase followed by an
517 other increase is more difficult to detect than an increase followed by a decrease) (Beaulieu et
518 al. 2008).

519 Here we focus on comparing EnvCpt to STARS and BMCpt, which have been used to
520 investigate changes in PDO and GMST, respectively. Overall, our approach clearly
521 outperforms these two methods. This result was to be expected as STARS and BMCpt only
522 consider a subset of the models fitted within EnvCpt. For example, the STARS methodology
523 is developed to detect shifts in the mean only. In terms of the model fit, it is equivalent to
524 considering only the Mean and Mean cpt models fitted with EnvCpt, thereby ignoring the
525 possibility of and misinterpreting underlying trends. BMCpt is more flexible than STARS
526 and designed to detect changes in the parameters of a regression model, so is also equivalent
527 to fitting the models Trend and Trend cpt. Since both of these approaches were developed for
528 independent data, all the models including an AR(1) are excluded from STARS and BMCpt.
529 While this issue can be mitigated with well-tuned pre-whitening (Appendix C), EnvCpt has
530 the additional advantage of natively supporting AR(1) detection without any parameter
531 tuning. In our attempts to tune the pre-whitening for STARS and BMCpt we used a sub-
532 sample size of 20, which is smaller than the length between the shifts inserted in the synthetic
533 series and shown to be optimal (Appendix C). Knowing *a priori* the minimum distance
534 between two shifts is of great benefit for the tuning, but the necessity of tuning is a great

535 disadvantage for STARS and BMCpt. That is, when the “truth” is unknown the choice of
536 parameter values for the pre-whitening is likely to induce errors (Fig. A2; Appendix C).

537 Several other methods have been proposed in the literature to detect multiple change-
538 points in environmental time-series (e.g. Beaulieu et al. 2012; Gazeaux et al. 2011; Lu et al.
539 2010; Reeves et al. 2007; Seidou and Ouarda 2007; Tomé and Miranda 2004; Wang 2008)
540 although these models assume independent errors and thus cannot distinguish signals from
541 autocorrelation, similar to STARS and BMCpt. To mitigate this issue one can use pre-
542 whitening techniques, although we show that pre-whitening has the disadvantage to
543 necessitate some parameters tuning. It has also been argued that an approach that forces the
544 lines of the piecewise linear model to meet assuring continuity between the trends is more
545 physically plausible in the case of the GMST (Cahill et al. 2015; Rahmstorf et al. 2017). Here,
546 we do not force the lines of the piecewise linear model to meet, but we find quasi-continuous
547 trends for the GMST (see Fig. 3). Imposing the continuity condition would restrain our
548 approach and make it unsuitable for the detection of climate regime shifts, which are
549 discontinuous and typically represented by abrupt changes in the mean. The main advantage
550 of the approach suggested here is its flexibility and applicability to a wide-range of climate
551 time-series, as illustrated through the GMST and PDO. The flexibility and breath of
552 applicability extends beyond inferring changes in the mean and trend as illustrated with these
553 two examples. Hence, EnvCpt is designed to detect change-points in all parameters of the
554 models fitted, including changes in autocorrelation and variance. There may be cases in
555 which the variability and/or dependence between successive observations are different after
556 the start of a new regime in the climate system or due to changes in measurements procedures.
557 Keeping the methodology as general as possible ensures these cases can also be analyzed
558 with EnvCpt.

559 Correctly identifying climate change signals is central to their understanding, as
560 mechanisms responsible for secular trends and abrupt changes are likely to be different (e.g.
561 anthropogenic influence vs natural forcings). However, abrupt changes can also be induced in
562 time-series through gradual increase in anthropogenic forcing when a critical threshold is
563 crossed (Lenton 2011). Further investigation of the forcing-response relationship can help
564 identify threshold and nonlinear dynamics, but correctly identifying the timing of an abrupt
565 change is a crucial first step (Andersen et al. 2009). Our EnvCpt approach is timely, as
566 increasing anthropogenic pressure on the climate system is expected to lead to more frequent
567 occurrences of abrupt changes in the physical climate system (Drijfhout et al. 2015).

568 Our methodology is flexible as it models different types of signals and memory in the
569 system. However, it assumes that temporal changes in climate time-series are piecewise
570 linear on a background of white noise or first-order autocorrelation, and that measurement
571 errors are random. While these assumptions are reasonable in many instances, there may be
572 cases of climate time-series with additional complexities such as long-term memory.
573 Departures from these assumptions may cause problems with the model selected as serious as
574 the five pervasive mistakes we are trying to avoid with EnvCpt. Thus, it is recommended to
575 combine the model selection with an analysis of the residuals as done here (Appendix E), and
576 to consider models that are physically plausible. Given that model selection is used with
577 EnvCpt, it can be easily extended to consider noise terms with additional parameters such as
578 an autoregressive moving-average (ARMA) models with higher-order and alternative model
579 forms (e.g. nonlinear). The models could be extended to take into account co-variables that
580 may explain part of the variability in climate time-series. For example, ENSO could
581 potentially explain part of the variability both in the GMST and PDO analyzed here, and
582 contribute to reducing the unexplained variability. When modifying the models used here,
583 one must keep in mind that the AIC weights are dependent on the subset of models being

584 compared. As such, if additional models were being considered, the probabilities of the eight
585 models compared here may change. Finally, another advantage of an approach based on
586 model selection is that it can be easily modified to use a different information criterion such
587 as the BIC, but the results may vary. We illustrate this in Appendix F and show that using the
588 BIC instead of the AIC in the simulation study can slightly improve the results for most cases
589 of synthetic series, except for the Mean $cpt + AR(1)$ case, for which the results are worst
590 (Figure A8). We refrain from making a universal recommendation here, as there are many
591 factors affecting the performance of AIC and BIC (Burnham and Anderson 2002) with
592 considerations that are going beyond our simulation study. This aspect should be the focus of
593 future work.

594

595 **Acknowledgements**

596 We thank SECURE and the EPSRC (EP/M008347/1) for funding this research. CB was also
597 supported by a Marie Curie FP7 Reintegration Grants within the Seventh European
598 Community Framework (project 631466 – TROPHYZ). The authors thank three anonymous
599 reviewers and E. Ruggieri for helpful comments that greatly improved the manuscript. We
600 thank S. Rodionov and E. Ruggieri for making their code freely available.

601 **APPENDIX A**

602 **Technical detail on the EnvCpt approach and simulations**

603 The EnvCpt approach fits eight different models to the data and returns the fit and
604 number of parameters for each model. The pseudo-code for the algorithm is as follows:

605 *EnvCpt Pseudo Algorithm*

606 Inputs: Time series y_t
607 msl = Minimum number of time points between changes (default 5)
608 pen = Penalty for changepoint algorithms (default MBIC)
609 Initialize: Let n = length of time series
610 Fit: 1. Constant mean with independent errors via maximum likelihood
611 2. Constant mean with AR(1) errors via maximum likelihood
612 3. Linear trend with independent errors via maximum likelihood
613 4. Linear trend with AR(1) errors via maximum likelihood
614 5. Constant mean changepoint model with independent errors via PELT
615 algorithm with msl and pen options.
616 6. Linear trend changepoint model with independent errors via PELT
617 algorithm with msl and pen options.
618 7. Constant mean changepoint model with AR(1) errors via PELT algorithm
619 with msl and pen options.
620 8. Linear trend changepoint model with AR(1) errors via PELT algorithm with
621 msl and pen options.
622 Output: A matrix of likelihood values and number of parameters for each model fit. A
623 list containing the fit for each of the eight models.

624 Using the output, one can compute an information criterion to determine the model that best
625 fits the data – in this study we use the AIC. See Appendix E for a sensitivity study to the
626 choice of criterion.

627 The PELT algorithm used in the EnvCpt procedure is described mathematically in
628 (Killick et al. 2012). Contrary to binary searches, where the most likely change is identified
629 and the time-series is split at that point, the PELT algorithm solves the segmentation problem
630 exactly by performing a search considering all options for any possible number of changes
631 (varying from 1 to the maximum number of change-points given the set minimum segment
632 length). This search is completed efficiently using a combination of dynamic programming
633 and pruning. Dynamic programming allows us to consider the data sequentially from the start
634 to the end and monitor the location of the last change-point only, which reduces the
635 computational time significantly. However, as the size of the data grows it remains time
636 consuming to monitor all potential last change-point locations. Thus, pruning is used to solve
637 this issue. For example, if there is an obvious change-point at, say time point 57, then the
638 probability of the last change being before that (e.g. time point 15) is zero. The definition of
639 “obvious” is controlled by the penalty parameter – a larger value means that a change has to
640 be larger to be considered “obvious”. If “obvious” changes occur throughout the data then
641 this dramatically reduces the computational time.

642 To evaluate the approach, we generate synthetic series from each one of the eight
643 models considered with parameters mimicking the GMST and PDO. For reproducibility, the
644 parameters used are presented in Table A1.

645

646 APPENDIX B

647 Choice of parameters for BMCpt

648 Hyper-parameters for the prior distributions of the regression parameters and variance
649 used with BMCpt are set following previous recommendations (Ruggieri 2012). We set the
650 variance scaling hyper-parameter for the multivariate Normal prior on the regression
651 parameters to 0.01. The hyper-parameters for the variance prior, i.e. the prior variance (σ_0^2),
652 is set to the variance of the data set being used. As for the pseudo data point of variance (ν_0),
653 which is recommended to be <25% of the minimum segment length (Ruggieri 2012), we vary
654 this parameter between 0 and 2.5 to find the value that optimizes the number of change-
655 points detected (Fig. A1). We focus on the number of change-points here, as these parameters
656 can affect the number of change-points detected, but not the distribution of their positions
657 (Ruggieri 2012). Tuning for ν_0 is performed for the four cases without AR(1) for which
658 BMCpt should perform well at identifying the true underlying model. For the cases scenario
659 with no change-points (i.e. Mean and Trend), the value of ν_0 does not have any impact on the
660 number of changes detected as none are detected for all values of ν_0 , thus these results are
661 not shown here. As illustrated in Fig. A1a, all values of ν_0 in the simulation scenario of a
662 trend with change-points (Trend cpt) lead to a low detection of the correct number of change-
663 points, but the most substantial improvement is obtained with a value of 0.25. In the case
664 scenario of mean change-points (Mean cpt), the correct number of change-points is obtained
665 at a highest frequency for any values of ν_0 (Fig. A1b). Setting ν_0 to 0 leads to no change-
666 points. Therefore, a value of 0.25 has been used subsequently in all simulations. Finally, the
667 maximum number of change-points is set to 10.

668 **APPENDIX C**

669 **Tuning of parameters for pre-whitening**

670 To reduce false alarms due to the presence of autocorrelation, pre-whitening of the
671 time-series was used with STARS and BMCpt (Rodionov 2006). This consists of removing
672 the first-order autocorrelation in the time-series such as:

673
$$x'_t = x_t - \hat{\rho}^c x_{t-1} \quad t = 2, \dots, n \quad (1)$$

674 where x_t and x'_t represent the raw and pre-whitened variable at time t respectively, n is the
675 length of the raw time-series and $\hat{\rho}^c$ represents the bias-corrected first-order autocorrelation
676 estimate. In a practical situation, the first-order autocorrelation used in pre-whitening is
677 unknown (and may also change over time). To obtain an estimate we used two approaches
678 developed by Marriott and Pope (1954) and Orcutt and Winokur Jr (1969), referred to as MP
679 and INV respectively. The MP estimate is given by:

680
$$\hat{\rho}^c = \frac{(m-1)\hat{\rho}+1}{(m-4)} \quad (2)$$

681 where $\hat{\rho}$ is the median of the first-order autocorrelation calculated in each subsample of size
682 m . The INV estimate uses four iterative corrections:

683
$$\hat{\rho}^{c,1} = \hat{\rho} + \frac{1}{m} \quad (3)$$

684
$$\hat{\rho}^{c,k} = \hat{\rho}^{c,k-1} + \frac{|\hat{\rho}^{c,k-1}|}{m} \quad k = 2,3,4 \quad (4)$$

685 In order to find an optimal value for the subsample size used in pre-whitening we conduct
686 simulations over a range of subsample sizes using the Mean cpt + AR(1) scenario. This is
687 done with both MP and INV approaches for pre-whitening using subsample sizes of 5, 10, 20,
688 30, 50 and 75 and illustrated in Figure A2. With both pre-whitening approaches, very large

689 (75) and very small (5) subsample size lead to a reduced rate of true positives and increased
690 false negatives towards the end of the time-series. A subsample size of approximately 20 is
691 shown optimal here, which is smaller than the distance between the two shifts (28 years).
692 When the number and location of changes is unknown, the choice of this parameter is rather
693 arbitrary and can have substantial effect on the results (Fig. A2).

694

695 **APPENDIX D**

696 **Results obtained after pre-whitening the synthetic data**

697 For comparison, we apply pre-whitening using both MP and INV in all simulations
698 with both STARS and BMCpt, and with a sub-sample size of 20, as chosen after optimization
699 (Fig. A2). Fig. A3 presents the number of shifts detected for the four simulation cases with
700 AR(1). For the two cases with no shifts: Trend + AR(1) and Mean + AR(1), BMCpt with pre-
701 whitening and EnvCpt are equivalent. The number of shifts detected is reduced for STARS,
702 but there is still a substantial rate of false detection. This is surprising, as STARS should be
703 able to return a no change model for the Mean + AR(1) case, but detects changes in over 34%
704 of the series. Nevertheless, the rate of false detection is reduced with pre-whitening, but
705 remains substantial with STARS. In presence of change-points (cases Trend cpt + AR(1) and
706 Mean cpt + AR(1)), the pre-whitening deteriorates BMCpt performance while it significantly
707 improves STARS ability to detect shifts in the mean.

708 Fig. A4 presents density estimates of the locations of the identified change-points for
709 synthetic series that were generated with change-points and AR(1). For the case Trend cpt +
710 AR(1), whilst the peaks of the true changes have a similar density to the EnvCpt method,
711 STARS and BMCpt tend to detect spurious changes towards the end of the series. In presence
712 of mean change-points, EnvCpt and both STARS and BMCpt applied with pre-whitening

713 succeed at identifying the correct timing of the change-points. While the densities in Fig. A4b
714 give the impression that BMCpt is performing better than STARS and EnvCpt with higher
715 peaks, this is due to fewer changes being detected with this approach (see Fig. A3d).

716

717 **APPENDIX E**

718 **Goodness-of-fit of the GMST and PDO best models**

719 To validate the models selected, we also verify their underlying assumptions of
720 normality and independence of the residuals with additional testing (Table A2). In all cases,
721 the normality assumption of the residuals is respected, but not the independence for all Trend
722 cpt fits on the GMST and the MLOST Trend cpt + AR(1) fits. To further investigate the
723 autocorrelation structure of the residuals for both the Trend cpt and Trend cpt +AR(1) fits,
724 the autocorrelation and partial autocorrelation functions are presented in Figs. A5-A6,
725 respectively. The autocorrelation and partial autocorrelation functions are consistent with the
726 tests of independence presented in Table A2: the residuals of the Trend cpt + AR(1) fits are
727 independent overall (except for the MLOST dataset) (Fig. A5), while the residuals of the
728 Trend cpt fit are not (Fig. A6). The autocorrelation and partial autocorrelation functions for
729 the HadCRUT4 and HadCRUT4krig datasets (Fig. A6a-b) reveals potential presence of a
730 second-order autocorrelation process (AR(2)) in the residuals. Therefore, our models were
731 also fitted with an AR(2) in the background such as : Mean + AR(2), Trend + AR(2), Mean
732 cpt + AR(2) and Trend cpt + AR(2). Table A3 presents the AIC differences of the models
733 fitted with a background AR(2) as opposed to the previously selected models (Trend cpt and
734 Trend cpt + AR(1); Table 1). These results show that despite a potential AR(2) structure in
735 the residuals, there is no benefit from adding an extra parameter to explain the autocorrelation
736 structure. The AIC differences for the models including an AR(2) are substantially larger

737 than those of the best models selected, i.e mostly larger than 10 indicating essentially no
738 evidence for choosing these models instead. There is one exception for the GISTEMP dataset,
739 for which the Trend cpt +AR(2) model has a Δ of 2.5, which suggests some evidence for this
740 model being the best, but not enough to be at play. Overall, for the five GMST datasets, the
741 Trend cpt + AR(1) fit provides the smallest AIC and meet the underlying assumptions of the
742 model. As for the PDO, the model with the smallest AIC (Mean + AR(1)) respects the
743 underlying assumptions of normality and independence (Fig. A7; Table A2).

744

745 **APPENDIX F**

746 **Sensitivity to the model selection criterion**

747 To evaluate the sensitivity to the choice of model selection criterion, we compare the
748 results obtained on all sets of synthetic series with EnvCpt using the Bayesian Information
749 Criterion (BIC) (Figure A8). In most cases, the EnvCpt performance is slightly improved
750 when using the BIC, except for the Mean cpt + AR(1) case for which the BIC detects no
751 change-points in strong majority while there are two.

752 We also calculate the BIC for the eight models fitted within EnvCpt to the GMST and
753 PDO datasets (Table A4), For all GMST datasets the model with the smallest BIC is Trend
754 cpt + AR(1). This result is slightly different than the results obtained using the AIC for the
755 HADCRUT4 dataset for which the Trend cpt model has the smallest AIC (Table 1). However,
756 we discarded the Trend cpt model for the HADCRUT4 dataset due to the presence of
757 autocorrelation in the residuals (Table A2; Figs A5-A6) and concluded that the second best
758 model, Trend cpt + AR(1), was more appropriate. Thus, the best models identified using the
759 BIC are consistent with the results obtained with the AIC (Figure 3).

760 **References**

- 761 Akaike, H., 1974: A new look at the statistical model identification. *IEEE Transactions on*
762 *Automatic Control*, **19**, 716-723.
- 763 Andersen, T., J. Carstensen, E. Hernandez-Garcia, and C. M. Duarte, 2009: Ecological
764 thresholds and regime shifts: approaches to identification. *Trends Ecol Evol*, **24**, 49-57.
- 765 Beaulieu, C., J. Chen, and J. L. Sarmiento, 2012: Change-point analysis as a tool to detect
766 abrupt climate variations. *Philos Trans A Math Phys Eng Sci*, **370**, 1228-1249.
- 767 Beaulieu, C., O. Seidou, T. B. M. J. Ouarda, X. Zhang, G. Boulet, and A. Yagouti, 2008:
768 Intercomparison of homogenization techniques for precipitation data. *Water Resources*
769 *Research*, **44**, W02425.
- 770 Beaulieu, C., and Coauthors, 2016: Marine regime shifts in ocean biogeochemical models: a
771 case study in the Gulf of Alaska. *Biogeosciences*, **13**, 4533-4553.
- 772 Burnham, K. P., and D. R. Anderson, 2002: *Model selection and multimodel inference.*
773 *Apractical information-theoretic approach*. 2 ed. Springer, New York. Boulton, C., and T. M.
774 Lenton, 2015: Slowing down of North Pacific variability and its implications for abrupt
775 ecosystem change. *Proc Natl Acad Sci U S A*, **112**, 11496-11501.
- 776 Cahill, N., S. Rahmstorf, and A. C. Parnell, 2015: Change points of global temperature.
777 *Environmental Research Letters*, **10**, 084002.
- 778 Chatfield, C., 2003: *The Analysis of Time Series: An Introduction*. 7 ed. Chapman and
779 Hall/CRC.
- 780 Cowtan, K., and R. G. Way, 2014: Coverage bias in the HadCRUT4 temperature series and
781 its impact on recent temperature trends. *Quarterly Journal of the Royal Meteorological*
782 *Society*, **140**, 1935-1944.
- 783 Drijfhout, S., and Coauthors, 2015: Catalogue of abrupt shifts in Intergovernmental Panel on
784 Climate Change climate models. *Proc Natl Acad Sci U S A*, **112**, E5777-5786.
- 785 Drijfhout, S. S., A. T. Blaker, S. A. Josey, A. J. G. Nurser, B. Sinha, and M. A. Balmaseda,
786 2014: Surface warming hiatus caused by increased heat uptake across multiple ocean basins.
787 *Geophysical Research Letters*, **41**, 7868-7874.
- 788 Faghmous, J. H., and V. Kumar, 2014: A Big Data Guide to Understanding Climate Change:
789 The Case for Theory-Guided Data Science. *Big Data*, **2**, 155-163.
- 790 Frankignoul, C., and K. Hasselmann, 1977: Stochastic climate models, Part II: Application to
791 sea-surface temperature anomalies and thermocline variability. *Tellus*, **29**, 289-305.

792 Franzke, C., 2012: Nonlinear trends, long-range dependence, and climate noise properties of
793 surface temperature. *Journal of Climate*, **25**, 4172-4183.

794 Fyfe, J. C., and Coauthors, 2016: Making sense of the early-2000s warming slowdown.
795 *Nature Climate Change*, **6**, 224-228.

796 Gazeaux, J., E. Flaounas, P. Naveau, and A. Hannart, 2011: Inferring change points and
797 nonlinear trends in multivariate time series: Application to West African monsoon onset
798 timings estimation. *Journal of Geophysical Research*, **116**, D05101.

799 Hansen, J., R. Ruedy, M. Sato, and K. Lo, 2010: *Global surface temperature change*. Vol. 48,
800 RG4004 pp.

801 Hartmann, D. L., and Coauthors, 2013: Observations: Atmosphere and Surface. *Climate*
802 *Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth*
803 *Assessment Report of the Intergovernmental Panel on Climate Change*, T. F. Stocker, and
804 Coauthors, Eds., Cambridge University Press.

805 Hasselmann, K., 1976: Stochastic climate models Part I. Theory. *Tellus*, **28**, 473-485.

806 Haynes, K., I. A. Eckley, and P. Fearnhead, 2017: Computationally Efficient Change-point
807 Detection for a Range of Penalties. *Journal of Computational and Graphical Statistics*, **26**,
808 134-143.

809 Huang, B., and Coauthors, 2015: Extended Reconstructed Sea Surface Temperature Version
810 4 (ERSST.v4). Part I: Upgrades and Intercomparisons. *Journal of Climate*, **28**, 911-930.

811 Huber, M., and R. Knutti, 2014: Natural variability, radiative forcing and climate response in
812 the recent hiatus reconciled. *Nature Geoscience*, **7**, 651-656.

813 Jones, G. S., and J. J. Kennedy, 2017: Sensitivity of Attribution of Anthropogenic Near-
814 Surface Warming to Observational Uncertainty. *Journal of Climate*, **30**, 4677-4691.

815 Jones, P., 2016: The reliability of global and hemispheric surface temperature records.
816 *Advances in Atmospheric Sciences*, **33**, 269-282.

817 Jones, P. D., D. H. Lister, T. J. Osborn, C. Harpham, M. Salmon, and C. P. Morice, 2012:
818 Hemispheric and large-scale land-surface air temperature variations: An extensive revision
819 and an update to 2010. *Journal of Geophysical Research: Atmospheres*, **117**, n/a-n/a.

820 Karl, T. R., K. R. W., and B. Baker, 2000: *The record breaking global temperatures of 1997*
821 *and 1998: Evidence for an increase in the rate of global warming?* Vol. 27, 719-722 pp.

822 Karl, T. R., and Coauthors, 2015: Possible artifacts of data biases in the recent global surface
823 warming hiatus. *Science*, **348**, 1469-1472.

824 Kellogg, W. W., 1993: An apparent moratorium on the greenhouse warming due to the deep
825 ocean. *Climatic Change*, **25**, 85-88.

826 Kennedy, J. J., N. A. Rayner, R. O. Smith, D. E. Parker, and M. Saunby, 2011a: Reassessing
827 biases and other uncertainties in sea surface temperature observations measured in situ since
828 1850: 2. Biases and homogenization. *Journal of Geophysical Research*, **116**.

829 ———, 2011b: Reassessing biases and other uncertainties in sea surface temperature
830 observations measured in situ since 1850: 1. Measurement and sampling uncertainties.
831 *Journal of Geophysical Research*, **116**.

832 Kent, E. C., and Coauthors, 2017: A Call for New Approaches to Quantifying Biases in
833 Observations of Sea Surface Temperature. *Bulletin of the American Meteorological Society*,
834 **98**, 1601-1616.

835 Killick, R., P. Fearnhead, and I. A. Eckley, 2012: Optimal detection of changepoints with a
836 linear computational cost. *Journal of the American Statistical Association*, **107**, 1590-1598.

837 Killick, R., C. Beaulieu, and S. Taylor, 2016: version 0.1.

838 Knutson, T. R., R. Zhang, and L. Horowitz, 2016: Prospects for a prolonged slowdown in
839 global warming in the early 21st century. *Nature Communications*, **7**, 13676.

840 Lean, J. L., and D. H. Rind, 2009: How will Earth's surface temperature change in future
841 decades? *Geophysical Research Letters*, **36**.

842 Lenton, T. M., 2011: Early warning of climate tipping points. *Nature Climate Change*, **1**,
843 201-209.

844 Lenton, T. M., V. Dakos, S. Bathiany, and M. Scheffer, 2017: Observed trends in the
845 magnitude and persistence of monthly temperature variability. *Scientific Reports*, **7**, 5940.

846 Lewandowsky, S., J. S. Risbey, and N. Oreskes, 2016: The “pause” in global warming:
847 turning a routine fluctuation into a problem for science. *Bulletin of the American*
848 *Meteorological Society*, **97**, 723-733.

849 Lewandowsky, S., N. Oreskes, J. S. Risbey, B. R. Newell, and M. Smithson, 2015: Seepage:
850 Climate change denial and its effect on the scientific community. *Global Environmental*
851 *Change*, **33**, 1-13.

852 Liu, W., and Coauthors, 2015: Extended Reconstructed Sea Surface Temperature Version 4
853 (ERSST.v4): Part II. Parametric and Structural Uncertainty Estimations. *Journal of Climate*,
854 **28**, 931-951.

855 Løvsletten, O., and Rypdal, M., 2016: Statistics of regional surface temperatures post year
856 1900. Long-range versus short-range dependence, and significance of warming trends.
857 *Journal of Climate*, **29**, 4057–4068, 2016.

858 Lu, Q., R. Lund, and T. C. M. Lee, 2010: An MDL approach to the climate segmentation
859 problem. *The Annals of Applied Statistics*, **4**, 299-319.

860 Lund, R., and J. Reeves, 2002: Detection of undocumented changepoints: a revision of the
861 two-phase model. *Journal of Climate*, **15**, 2547-2554.

862 Mantua, N. R., S. R. Hare, Y. Zhang, J. M. Wallace, and R. C. Francis, 1997: A Pacific
863 interdecadal oscillation with impacts on salmon production. *Bulletin of the American*
864 *Meteorological Society*, **78**, 1069-1079.

865 Marriott, F. H. C., and J. A. Pope, 1954: Bias in the estimation of autocorrelations.
866 *Biometrika*, **41**, 390-402.

867 Medhaug, I., M. B. Stolpe, E. M. Fischer, and R. Knutti, 2017: Reconciling controversies
868 about the global warming hiatus'. *Nature*, **545**, 41-47.

869 Meehl, G. A., H. Teng, and J. M. Arblaster, 2014: Climate model simulations of the observed
870 early-2000s hiatus of global warming. *Nature Climate Change*, **4**, 898-902.

871 Morice, C. P., J. J. Kennedy, N. A. Rayner, and P. D. Jones, 2012: Quantifying uncertainties
872 in global and regional temperature change using an ensemble of observational estimates: The
873 HadCRUT4 data set. *Journal of Geophysical Research: Atmospheres*, **117**, D08101.

874 Mustin, K., C. Dytham, T. G. Benton, J. M. J. Travis, and J. Watson, 2013: Red noise
875 increases extinction risk during rapid climate change. *Diversity and Distributions*, **19**, 815-
876 824.

877 Newman, M., and Coauthors, 2016: The Pacific Decadal Oscillation, revisited. *Journal of*
878 *Climate*, **29**, 4399-4427.

879 NRC, 2013: *Abrupt impacts of climate change: anticipating surprises*. The National
880 Academies Press.

881 Orcutt, G. H., and H. S. Winokur Jr, 1969: First order autoregression: inference, estimation
882 and prediction. *Econometrica*, **37**, 1-14.

883 Poppick, A., E. J. Moyer, and M. L. Stein, 2017: Estimating trends in the global mean
884 temperature record. *Advances in Statistical Climatology, Meteorology and Oceanography*, **3**,
885 33-53.

886 Rahmstorf, S., G. Foster, and N. Cahill, 2017: Global temperature evolution: recent trends
887 and some pitfalls. *Environmental Research Letters*, **12**, 054001.

888 Rajaratnam, B., J. Romano, M. Tsiang, and N. S. Diffenbaugh, 2015: Debunking the climate
889 hiatus. *Climatic Change*, **133**, 129-140.

890 Reeves, J., J. Chen, X. L. Wang, R. Lund, and Q. Lu, 2007: A review and comparison of
891 changepoint detection techniques for climate data. *Journal of Applied Meteorology and*
892 *Climatology*, **46**, 900-915.

893 Risbey, J. S., S. Lewandosky, C. Langlais, D. P. Monselesan, T. J. O’Kane, and N. Oreskes,
894 2014: Well-estimated global surface warming in climate projections selected for ENSO phase.
895 *Nature Climate Change*, **4**, 835-840.

896 Robbins, M. W., C. M. Gallagher, and R. B. Lund, 2016: A general regression changepoint
897 test for time series data. *Journal of the American Statistical Association*, **111**, 670-683.

898 Rodionov, S. N., 2004: A sequential algorithm for testing climate regime shifts. *Geophysical*
899 *Research Letters*, **31**, L09204.

900 ———, 2006: Use of prewhitening in climate regime shift detection. *Geophysical Research*
901 *Letters*, **33**.

902 Rohde, R., and Coauthors, 2013: A New Estimate of the Average Earth Surface Land
903 Temperature Spanning 1753 to 2011. *Geoinformatics & Geostatistics: An Overview*, **01**.

904 Rudnick, D. L., and R. E. Davis, 2003: Red noise and regime shifts. *Deep-Sea Research Part*
905 *I*, **50**, 691-699.

906 Ruggieri, E., 2012: A Bayesian approach to detecting change points in climatic records.
907 *International Journal of Climatology*, **33**, 520-528.

908 Santer, B. D., and Coauthors, 2014: Volcanic contribution to decadal changes in tropospheric
909 temperature. *Nature Geoscience*, **7**, 185-189.

910 Schmidt, G. A., D. T. Shindell, and K. Tsigaridis, 2014: Reconciling warming trends. *Nature*
911 *Geoscience*, **7**, 158-160.

912 Schwarz, G., 1978: Estimating the dimension of a model. *The Annals of Statistics*, **6**, 461-
913 464.

914 Seidel, D. J., and J. R. Lanzante, 2004: An assessment of three alternatives to linear trends for
915 characterizing global atmospheric temperature changes. *Journal of Geophysical Research*,
916 **109**, D14108.

917 Seidou, O., and T. B. M. J. Ouarda, 2007: Recursion-based multiple changepoint detection in
918 multiple linear regression and application to river streamflows. *Water Resources Research*,
919 **43**, W07404.

920 Serinaldi, F., and C. G. Kilsby, 2016: The importance of prewhitening in change point
921 analysis under persistence. *Stochastic Environmental Research and Risk Assessment*, **30**, 763-
922 777.

923 Smith, T. M., R. W. Reynolds, T. R. Peterson, and J. H. Lawrimore, 2008: Improvements to
924 NOAA’s historical merged and–ocean surface temperature analysis (1880–2006). *Journal of*
925 *Climate*, **21**, 2283-2296.

926 Thompson, D. W., J. J. Kennedy, J. M. Wallace, and P. D. Jones, 2008: A large discontinuity
927 in the mid-twentieth century in observed global-mean surface temperature. *Nature*, **453**, 646-
928 649.

929 Tomé, A. R., and P. M. A. Miranda, 2004: Piecewise linear fitting and trend changing points
930 of climate parameters. *Geophysical Research Letters*, **31**, L02207.

931 Trenberth, K. E., 2015: Has there been a hiatus? *Science*, **349**, 691-692.

932 Trenberth, K. E., and J. T. Fasullo, 2013: An apparent hiatus in global warming? *Earth's*
933 *Future*, **1**, 19-32.

934 Vallis, G. K., 2010: Mechanisms of climate variability from years to decade. *Stochastic*
935 *Physics and Climate Modelling*, T. Palmer, and P. Williams, Eds., Cambridge University
936 Press, 496.

937 von Storch, H., 1999: Misuses of statistical analysis in climate research. *Analysis of Climate*
938 *Variability*, H. von Storch, and A. Navarra, Eds., Springer, 11-26.

939 von Storch, H., and F. W. Zwiers, 1999: *Statistical analysis in climate research*. Cambridge
940 University Press, 455 pp.

941 Vose, R. S., and Coauthors, 2012: NOAA's merged land-ocean surface temperature analysis.
942 *Bulletin of the American Meteorological Society*, 1677-1685.

943 Wang, S., J. Huang, Y. He, and Y. Guan, 2014: Combined effects of the Pacific Decadal
944 Oscillation and El Nino-Southern Oscillation on global land dry-wet changes. *Sci Rep*, **4**,
945 6651.

946 Wang, X. L., 2008: Accounting for Autocorrelation in Detecting Mean Shifts in Climate Data
947 Series Using the Penalized Maximal t-Test. *Journal of Applied Meteorology and*
948 *Climatology*, **47**, 2423-2444.

949 Wang, X. L., Q. H. Wen, and Y. Wu, 2007: Penalized maximal t test for detecting
950 undocumented mean change in climate data series. *Journal of Applied Meteorology and*
951 *Climatology*, **46**, 916-931.

952 Wang, X. L., H. Chen, Y. Wu, Y. Feng, and Q. Pu, 2010: New Techniques for the Detection
953 and Adjustment of Shifts in Daily Precipitation Data Series. *Journal of Applied Meteorology*
954 *and Climatology*, **49**, 2416-2436.

955 Wunsch, C., 1999: The interpretation of short climate records, with comments on the North
956 Atlantic and Southern Oscillations. *Bulletin of the American Meteorological Society*, **80**, 245-
957 255.

- 958 Yuan, N., Ding, M., Huang, Y., Fu, Z., Xoplaki, E., and J. Luterbacher, 2015: On the Long-
959 Term Climate Memory in the Surface Air Temperature Records over Antarctica: A
960 Nonnegligible Factor for Trend Evaluation. *Journal of Climate*, **28**, 5922–5934.
- 961 Zhang, N. R., and D. O. Siegmund, 2007: A Modified Bayes Information Criterion with
962 Applications to the Analysis of Comparative Genomic Hybridization Data. *Biometrics*, **63**,
963 22-32.
- 964 Zhang, Y., J. M. Wallace, and D. S. Battisti, 1997: ENSO-like Interdecadal Variability:
965 1900–93. *Journal of Climate*, **10**, 1004-1020.

966 **Tables**

967 Table 1: Comparison of the eight EnvCpt models on the GMST and PDO datasets. AIC
 968 differences (Δ) between the model with the smallest AIC and the seven other models, as well
 969 as their Akaike weights (w) representing the probabilities of each model being the best model
 970 given the data and the set of models considered. The model with the smallest AIC has a Δ of 0
 971 and is indicated in bold along with its associated probability. Blanks are left for change-point
 972 models that did not detect change-points, as the model fit is the same as the equivalent model
 973 without change-points.

Model	Data											
	HadCRUT4		HadCRUT4krig		BEST		MLOST		GISTEMP		PDO	
	Δ	w	Δ	w	Δ	w	Δ	w	Δ	w	Δ	w
1.Mean	355.5	0.00	372.7	0.00	386.5	0.00	340.6	0.00	326.7	0.00	42.5	0.00
2.Mean + AR(1)	46.0	0.00	40.7	0.00	40.0	0.00	35.8	0.00	38.5	0.00	0.0	0.56
3.Trend	165.2	0.00	162.2	0.00	150.3	0.00	152.1	0.00	136.9	0.00	44.5	0.00
4.Trend +AR(1)	31.3	0.00	25.9	0.00	23.3	0.00	23.2	0.00	24.6	0.00	1.1	0.44
5.Mean cpt	40.7	0.00	45.7	0.00	25.3	0.00	61.3	0.00	43.2	0.00	25.8	0.00
6.Mean cpt +AR(1)												
7.Trend cpt	0.0	0.98	1.5	0.32	16.8	0.00	26.0	0.00	13.4	0.00	23.4	0.00
8.Trend cpt +AR(1)	7.8	0.02	0.0	0.68	0.0	1.00	0.0	1.00	0.0	1.00		

975
 976

977 Table 2: Trend and first-order autocorrelation (AR(1)) parameter estimates for the model
 978 with trend change-points and AR(1) (Trend cpt + AR(1)) in the five GMST datasets.

Dataset	Cpt timing	Trend		AR(1)	
		Before cpt	After cpt	Before cpt	After cpt
HadCRUT4	1962	0.001	0.013	0.653	0.195
HadCRUT4krig	1972	0.001	0.018	0.635	0.083
BEST	1962	0.001	0.015	0.656	0.148
MLOST	1962	0.001	0.015	0.706	0.144
GISTEMP	1962	0.002	0.016	0.644	0.112

979

980 Table A1: List of parameters used to simulate the sets of synthetic series.

981

Variable	Model	Parameters
PDO (n=116 years)	Mean	$\mu = 0.028, \sigma = 0.8$
	Mean + AR(1)	$\mu = 0.049, \varphi = 0.522, \sigma = 0.8$
	Mean cpt	$\mu_1 = 0.222, \mu_2 = -0.652, \mu_3 = 0.271$ $c_1 = 49, c_2 = 77, m = 3, \sigma = 0.3$
	Mean cpt + AR(1)	$\mu_1 = 0.222, \mu_2 = -0.652, \mu_3 = 0.271$ $\varphi_1 = \varphi_2 = \varphi_3 = 0.402$ $c_1 = 49, c_2 = 77, m = 3, \sigma = 0.3$
GMST (n=166 years)	Trend	$\lambda = -0.513, \beta = 0.005, \sigma = 0.1$
	Trend + AR(1)	$\lambda = -0.128, \beta = 0.001, \varphi = 0.756, \sigma = 0.3$
	Trend cpt	$\lambda_1 = -0.299, \lambda_2 = -1.327, \lambda_3 = 0.171, \lambda_4 = -2.124,$ $\beta_1 = -0.001, \beta_2 = 0.014, \beta_3 = -0.002,$ $\beta_4 = 0.016, c_1 = 57, c_2 = 96, c_3 = 127, m = 4,$ $\sigma = 0.4$
	Trend cpt + AR(1)	$\lambda_1 = -0.112, \lambda_2 = -1.707, \beta_1 = -0.001,$ $\beta_2 = 0.013, \varphi_1 = 0.659, \varphi_2 = 0.153, c_1 = 113,$ $m = 2, \sigma = 0.1$

982

983 Table A2: Results (p-value) of the Lilliefors (L) and Durbin-Watson (DW) tests applied to
 984 the residuals of the best models fitted to the GMST (Trend cpt and Trend cpt + AR(1) and
 985 PDO datasets (Mean + AR(1)).

Model	Test	Data					
		HadCRUT4	HadCRUT4krig	BEST	MLOST	GISTEMP	PDO
Trend cpt	L	0.50	0.50	0.29	0.39	0.12	
	DW	<0.001*	<0.001*	<0.001*	<0.001*	<0.001*	
Trend cpt + AR(1)	L	0.39	0.50	0.33	0.50	0.08	
	DW	0.53	0.25	0.19	<0.001*	0.66	
Mean + AR(1)	L						0.50
	DW						0.68

986 *Significant at the 1% critical level.

987 Table A3: Comparison of the best EnvCpt models (Trend cpt and Trend cpt + AR(1)) with
 988 models including a second-order autocorrelation process (AR(2)) on the GMST and PDO
 989 datasets. AIC differences (Δ) between the model with the smallest AIC and the other models
 990 are presented. The model with the smallest AIC has a Δ of 0 and is indicated in bold.

991

Model	Data				
	HadCRUT4	HadCRUT4krig	BEST	MLOST	GISTEMP
Trend cpt	0.0	1.5	16.8	26.0	13.5
Trend cpt + AR(1)	7.8	0.0	0.0	0.0	0.0
Mean + AR(2)	41.6	37.1	37.5	34.4	35.5
Trend + AR(2)	30.5	25.0	24.8	25.4	25.2
Mean cpt + AR(2)	48.0	47.7	42.1	37.8	40.5
Trend cpt + AR(2)	42.5	37.0	36.8	37.4	2.5

992

993 Table A4: Bayesian Information Criterion (BIC) differences for the eight models within
 994 EnvCpt fitted to the GMST and PDO datasets. The model with the smallest BIC has a Δ of 0
 995 and is indicated in bold. Blanks are left for change-point models that did not detect change-
 996 points, as the model fit is the same as the equivalent model without change-points.

Model	Data					
	HadCRUT4	HadCRUT4krig	BEST	MLOST	GISTEMP	PDO
1. Mean	325.8	350.9	364.7	320.2	307.4	39.1
2. Mean + AR(1)	19.5	22.0	21.3	18.3	24.1	0.0
3. Trend	138.6	143.6	131.6	134.6	-122.6	43.9
4. Trend +AR(1)	7.8	10.3	7.7	8.6	13.0	3.3
5. Mean cpt	39.1	51.9	40.9	67.1	44.8	30.7
6. Mean cpt +AR(1)						
7. Trend cpt	10.8	20.2	23.0	51.5	23.3	33.8
8. Trend cpt +AR(1)	0.0	0.0	0.0	0.0	0.0	

998

999 **Figure Captions**

1000

1001 **Figure 1:** Five possible misuses of statistics when inferring changes in climate time-series
1002 exhibiting a long-term linear trend, shifts or memory: a) fitting a linear trend in presence of
1003 shifts in the mean or shifts in trend; b) fitting shifts in the mean in presence of a trend; c)
1004 fitting a linear trend assuming independent errors (i.e. white noise) in presence of
1005 autocorrelation; d) fitting shifts in the mean assuming white noise in presence of
1006 autocorrelation; e) fitting a first-order autocorrelation model in presence of mean shifts.

1007 **Figure 2:** Datasets used in this study a) global mean surface temperature (GMST) from the
1008 Met Office Hadley Centre surface temperature (HadCRUT4), HadCRUT4 infilled by kriging
1009 (HadCRUT4krig), Berkeley Earth Surface Temperature (BEST), Merged Land–Ocean
1010 Surface Temperature Analysis (MLOST), and Goddard Institute of Space Studies Surface
1011 Temperature Analysis (GISTEMP) and b) the Pacific Decadal Oscillation (PDO).

1012 **Figure 3:** Fit of the eight models in EnvCpt to five global mean surface temperature (GMST)
1013 datasets: a) Met Office Hadley Centre surface temperature (HadCRUT4), b) HadCRUT4
1014 infilled by kriging (HadCRUT4krig), c) Berkeley Earth Surface Temperature (BEST), d)
1015 Merged Land–Ocean Surface Temperature Analysis (MLOST), e) Goddard Institute of Space
1016 Studies Surface Temperature Analysis (GISTEMP) and f) the Pacific Decadal Oscillation
1017 (PDO). The tick marks indicate where change-points were detected. For each dataset, the
1018 Akaike Information Criterion differences (Δ) between each model and the best model
1019 (smallest AIC) are also shown on a logarithmic scale adjusted so that the best model has a log
1020 difference of zero, and is indicated by a star. The dotted vertical lines indicate cutoffs of
1021 models evidence: there is substantial support for models with a difference below the red line
1022 and essentially no support for models with differences above the black line.

1023 **Figure 4:** Synthetic time-series example from each simulation scenario case a) a linear trend,
1024 b) a linear trend with first-order autocorrelation, c) a trend with three change-points in the
1025 regression parameters, d) a trend with a change-point in the regression parameters and first-
1026 order autocorrelation, e) a constant mean, f) a constant mean with first-order autocorrelation,
1027 g) two change-points in the mean and h) two change-points in the mean with first-order
1028 autocorrelation. For each case, a total number of 1,000 random replications are simulated.

1029 **Figure 5:** Number of change-points detected with EnvCpt, STARS and BMCpt for each
1030 simulated scenario across 1,000 replications a) a linear trend, b) a linear trend with first-order
1031 autocorrelation, c) a trend with three change-points in the regression parameters, d) a trend
1032 with a change-point in the regression parameters and first-order autocorrelation, e) a constant
1033 mean, f) a constant mean with first-order autocorrelation, g) two change-points in the mean
1034 and h) two change-points in the mean with first-order autocorrelation. Overall, EnvCpt is
1035 closer to the true number of change-points than STARS and BMCpt.

1036 **Figure 6:** Density of change-point timings detected using EnvCpt, STARS and BMCpt for
1037 the four simulated scenarios with change-points across 1,000 replications a) a trend with
1038 three change-points in the regression parameters, b) a trend with a change-point in the
1039 regression parameters and first-order autocorrelation, c) two change-points in the mean and
1040 d) two change-points in the mean with first-order autocorrelation. Overall, EnvCpt identifies
1041 correctly the true change-point locations while STARS and BMCpt may detect change-points
1042 at timings when none were introduced in the synthetic series in presence of trend change-
1043 points.

1044 **Figure A1:** Number of change-points detected with BMCpt for the a) Trend cpt and b) Mean
1045 cpt scenario across 1,000 replications. Change-points were detected using a range of values
1046 for the pseudo data point of variance parameter (v_0). A value of 0.25 is shown optimal here.

1047 **Figure A2:** Density of change-point locations for the change-points in the mean and a
1048 background AR(1) (Mean cpt + AR(1)) scenario across 1,000 replications. Change-points
1049 were detected with a) STARS and b) BMCpt methodologies using a range of subsample sizes
1050 for pre-whitening using the MP and INV approaches. A subsample size of 20 is shown
1051 optimal here for both methods. For STARS, very large or very small subsample sizes lead to
1052 false detections at the end of the time-series. For BMCpt, very large or very small sample
1053 sizes lead to improved detection of one shift to the detriment of the other.

1054 **Figure A3:** Number of change-points detected with EnvCpt, and STARS and BMCpt with
1055 pre-whitening for each simulated scenario across 1,000 replications a) a linear trend, b) a
1056 linear trend with first-order autocorrelation, c) a trend with three change-points in the
1057 regression parameters, d) a trend with a change-point in the regression parameters and first-
1058 order autocorrelation, e) a constant mean, f) a constant mean with first-order autocorrelation,
1059 g) two change-points in the mean and h) two change-points in the mean with first-order
1060 autocorrelation. The pre-whitening is performed using the using the MP and INV approaches
1061 with a subsample size of 20.

1062 **Figure A4:** Density of change-point timings detected using EnvCpt, STARS and BMCpt
1063 with pre-whitening for the two simulated scenarios with change-points and AR(1) across
1064 1,000 replications a) a trend with a change-point in the regression parameters and first-order
1065 autocorrelation and b) two change-points in the mean with first-order autocorrelation. The
1066 pre-whitening is performed using the using the MP and INV approaches with a subsample
1067 size of 20.

1068 **Figure A5:** Autocorrelation and partial autocorrelation function of the residuals from the
1069 Trend cpt + AR(1) model fitted to the global mean surface temperature datasets a)
1070 HadCRUT4, b) HadCRUT4krig, c) BEST, d) MLOST and e) GISTEMP. Dashed lines

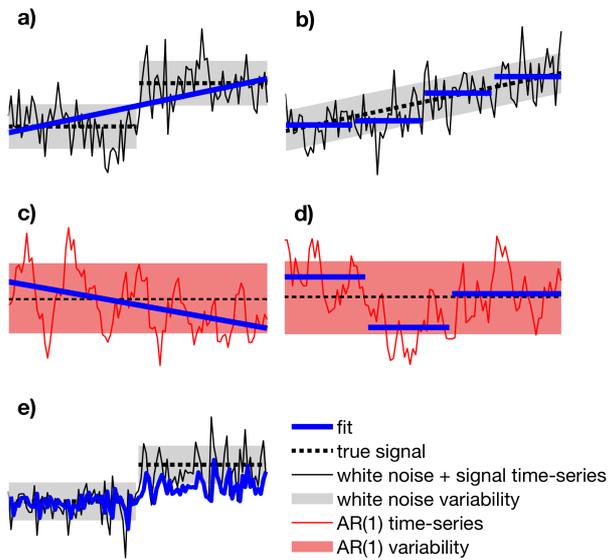
1071 represent the 95% confidence intervals on the partial autocorrelation.

1072 **Figure A6:** Autocorrelation and partial autocorrelation function of the residuals from the
1073 Trend cpt model fitted to the global mean surface temperature datasets a) HadCRUT4, b)
1074 HadCRUT4krig, c) BEST, d) MLOST and e) GISTEMP. Dashed lines represent the 95%
1075 confidence intervals on the partial autocorrelation.

1076 **Figure A7:** Autocorrelation and partial autocorrelation function of the residuals from the
1077 Mean + AR(1) model fitted to the PDO. Dashed lines represent the 95% confidence intervals
1078 on the partial autocorrelation.

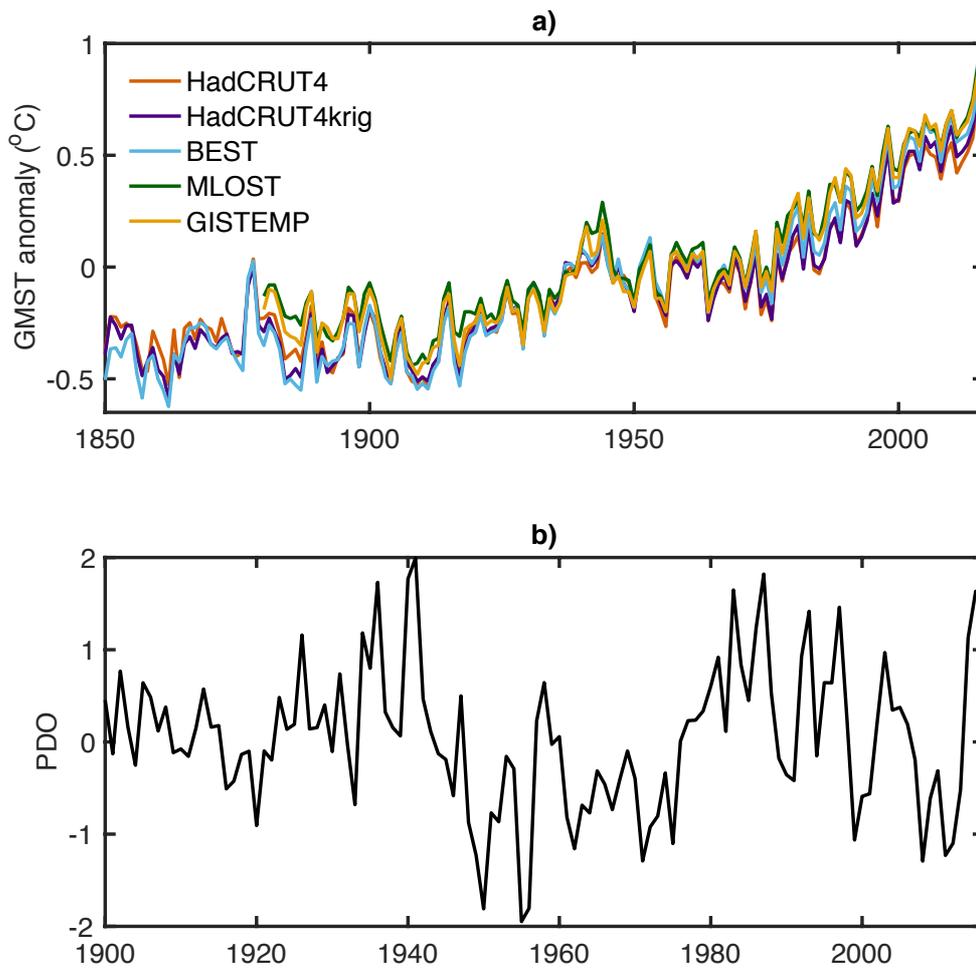
1079 **Figure A8:** Number of change-points detected with EnvCpt with either the Akaike
1080 Information Criterion (AIC) vs the Bayesian Information Criterion (BIC) for each simulated
1081 scenario across 1,000 replications a) a linear trend, b) a linear trend with first-order
1082 autocorrelation, c) a trend with three change-points in the regression parameters, d) a trend
1083 with a change-point in the regression parameters and first-order autocorrelation, e) a constant
1084 mean, f) a constant mean with first-order autocorrelation, g) two change-points in the mean
1085 and h) two change-points in the mean with first-order autocorrelation.

1086 **Figures**

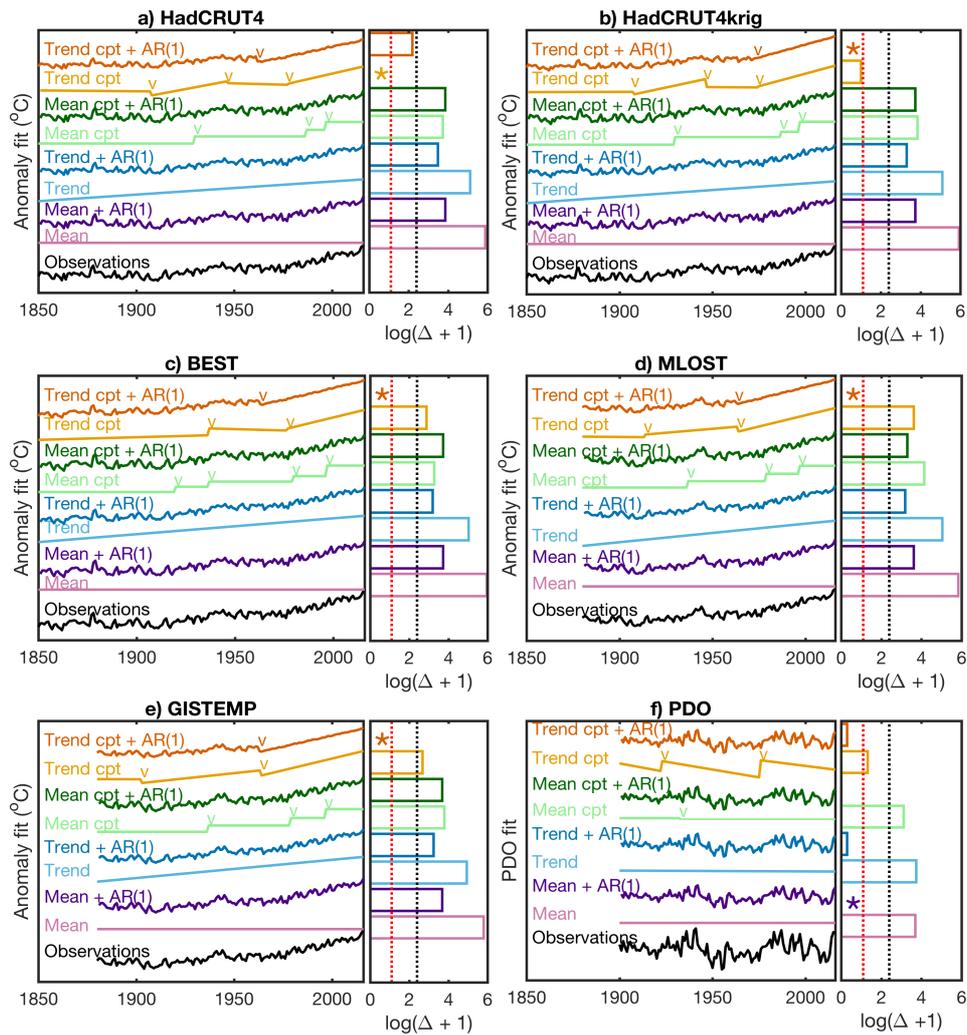


1087

1088 **Figure 1:** Five possible misuses of statistics when inferring changes in climate time-series
1089 exhibiting a long-term linear trend, shifts or memory: a) fitting a linear trend in presence of
1090 shifts in the mean or shifts in trend; b) fitting shifts in the mean in presence of a trend; c)
1091 fitting a linear trend assuming independent errors (i.e. white noise) in presence of
1092 autocorrelation; d) fitting shifts in the mean assuming white noise in presence of
1093 autocorrelation; e) fitting a first-order autocorrelation model in presence of mean shifts.



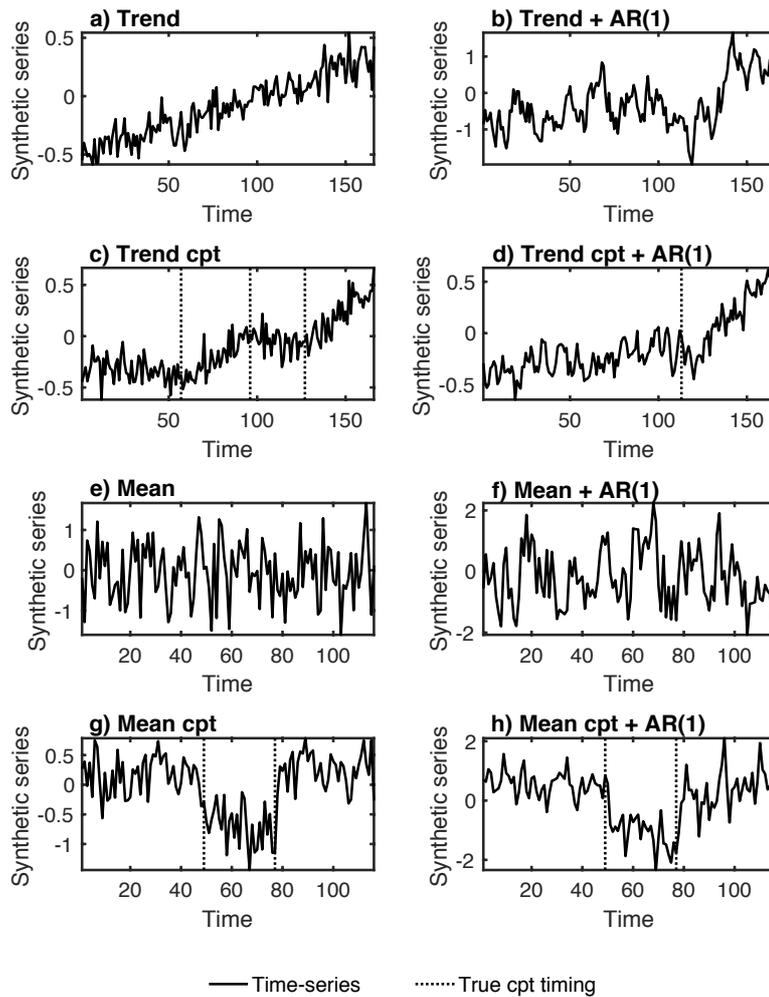
1094
 1095 **Figure 2:** Datasets used in this study a) global mean surface temperature (GMST) from the
 1096 Met Office Hadley Centre surface temperature (HadCRUT4), HadCRUT4 infilled by kriging
 1097 (HadCRUT4krig), Berkeley Earth Surface Temperature (BEST), Merged Land–Ocean
 1098 Surface Temperature Analysis (MLOST), and Goddard Institute of Space Studies Surface
 1099 Temperature Analysis (GISTEMP) and b) the Pacific Decadal Oscillation (PDO).



1100

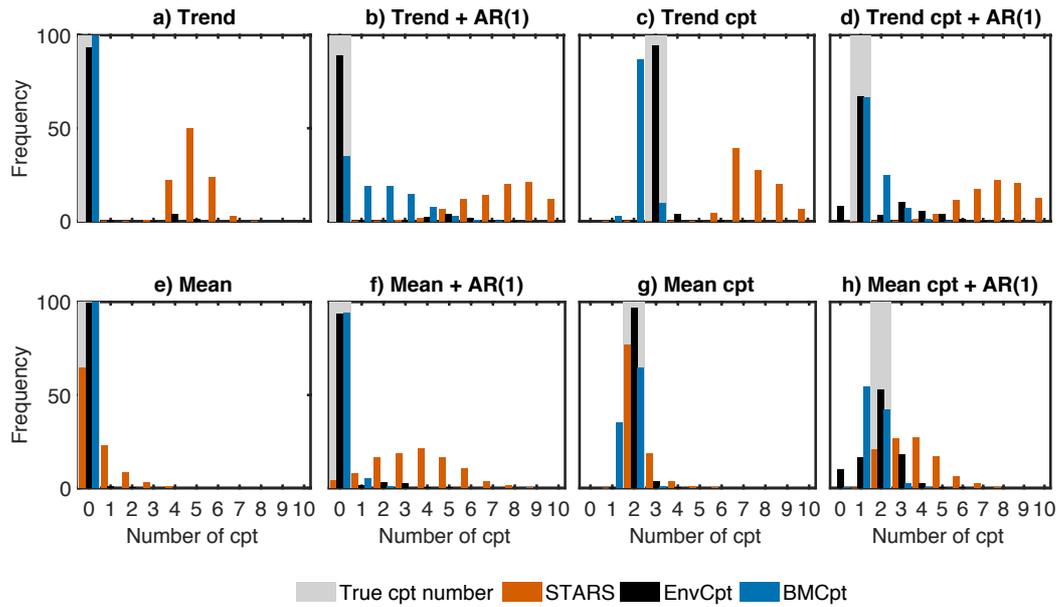
1101 **Figure 3:** Fit of the eight models in EnvCpt to five global mean surface temperature (GMST)
 1102 datasets: a) Met Office Hadley Centre surface temperature (HadCRUT4), b) HadCRUT4
 1103 infilled by kriging (HadCRUT4krig), c) Berkeley Earth Surface Temperature (BEST), d)
 1104 Merged Land–Ocean Surface Temperature Analysis (MLOST), e) Goddard Institute of Space
 1105 Studies Surface Temperature Analysis (GISTEMP) and f) the Pacific Decadal Oscillation
 1106 (PDO). The tick marks indicate where change-points were detected. For each dataset, the
 1107 Akaike Information Criterion differences (Δ) between each model and the best model
 1108 (smallest AIC) are also shown on a logarithmic scale adjusted so that the best model has a log
 1109 difference of zero, and is indicated by a star. The dotted vertical lines indicate cutoffs of

- 1110 models evidence: there is substantial support for models with a difference below the red line
- 1111 and essentially no support for models with differences above the black line.



1112

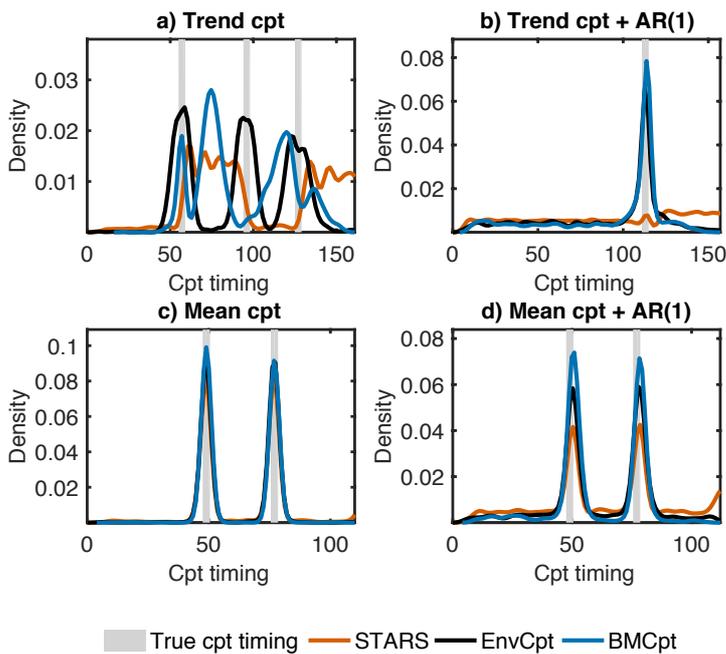
1113 **Figure 4:** Synthetic time-series example from each simulation scenario case a) a linear trend,
 1114 b) a linear trend with first-order autocorrelation, c) a trend with three change-points in the
 1115 regression parameters, d) a trend with a change-point in the regression parameters and first-
 1116 order autocorrelation, e) a constant mean, f) a constant mean with first-order autocorrelation,
 1117 g) two change-points in the mean and h) two change-points in the mean with first-order
 1118 autocorrelation. For each case, a total number of 1,000 random replications are simulated.



1119

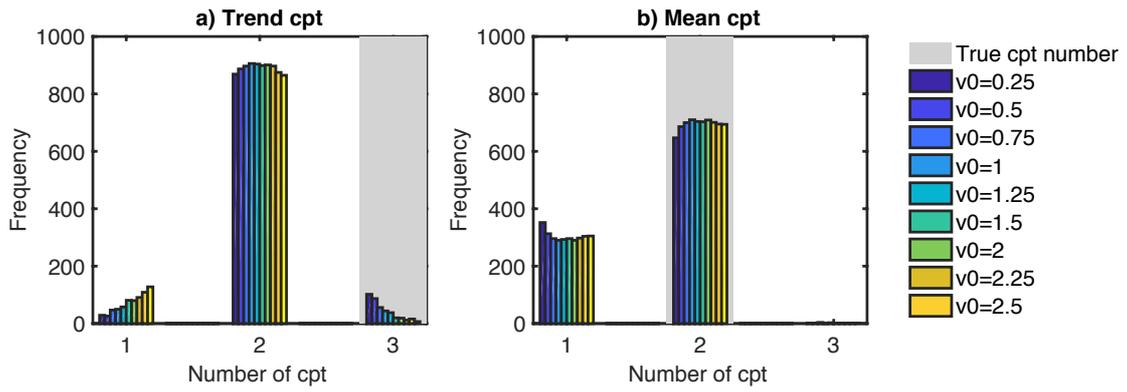
1120 **Figure 5:** Number of change-points detected with EnvCpt, STARS and BMCpt for each
 1121 simulated scenario across 1,000 replications a) a linear trend, b) a linear trend with first-order
 1122 autocorrelation, c) a trend with three change-points in the regression parameters, d) a trend
 1123 with a change-point in the regression parameters and first-order autocorrelation, e) a constant
 1124 mean, f) a constant mean with first-order autocorrelation, g) two change-points in the mean
 1125 and h) two change-points in the mean with first-order autocorrelation. Overall, EnvCpt is
 1126 closer to the true number of change-points than STARS and BMCpt.

1127



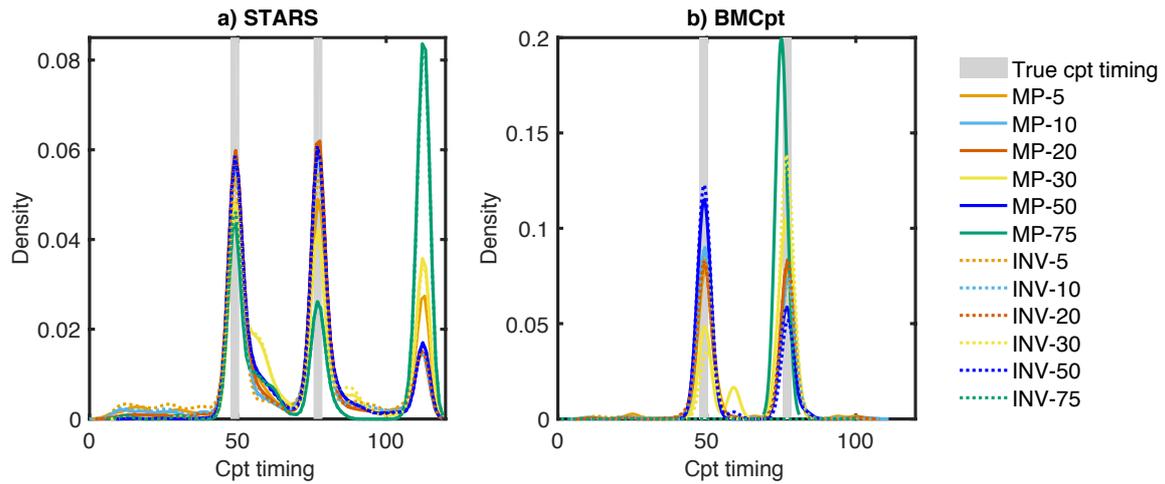
1128

1129 **Figure 6:** Density of change-point timings detected using EnvCpt, STARS and BMCpt for
 1130 the four simulated scenarios with change-points across 1,000 replications a) a trend with
 1131 three change-points in the regression parameters, b) a trend with a change-point in the
 1132 regression parameters and first-order autocorrelation, c) two change-points in the mean and
 1133 d) two change-points in the mean with first-order autocorrelation. Overall, EnvCpt identifies
 1134 correctly the true change-point locations while STARS and BMCpt may detect change-points
 1135 at timings when none were introduced in the synthetic series in presence of trend change-
 1136 points.



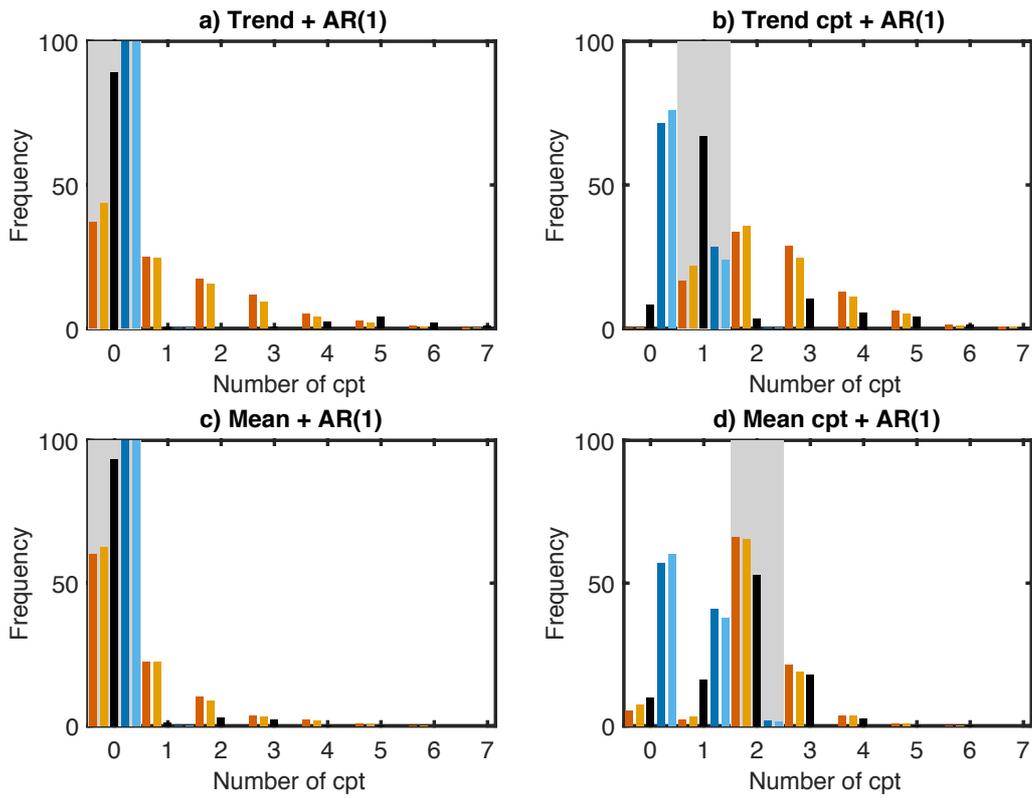
1137

1138 **Figure A1:** Number of change-points detected with BMCpt for the a) Trend cpt and b) Mean
 1139 cpt scenario across 1,000 replications. Change-points were detected using a range of values
 1140 for the pseudo data point of variance parameter (v_0). A value of 0.25 is shown optimal here.



1141

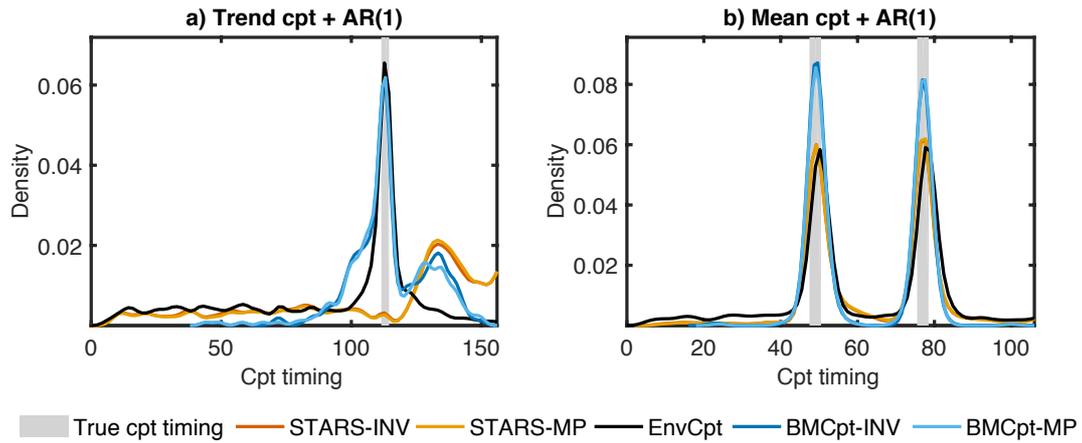
1142 **Figure A2:** Density of change-point locations for the change-points in the mean and a
 1143 background AR(1) (Mean cpt + AR(1)) scenario across 1,000 replications. Change-points
 1144 were detected with a) STARS and b) BMCpt methodologies using a range of subsample sizes
 1145 for pre-whitening using the MP and INV approaches. A subsample size of 20 is shown
 1146 optimal here for both methods. For STARS, very large or very small subsample sizes lead to
 1147 false detections at the end of the time-series. For BMCpt, very large or very small sample
 1148 sizes lead to improved detection of one shift to the detriment of the other.



Legend: True cpt number (grey shaded area), STARS-INV (orange), STARS-MP (yellow), EnvCpt (black), BMCpt-INV (dark blue), BMCpt-MP (light blue)

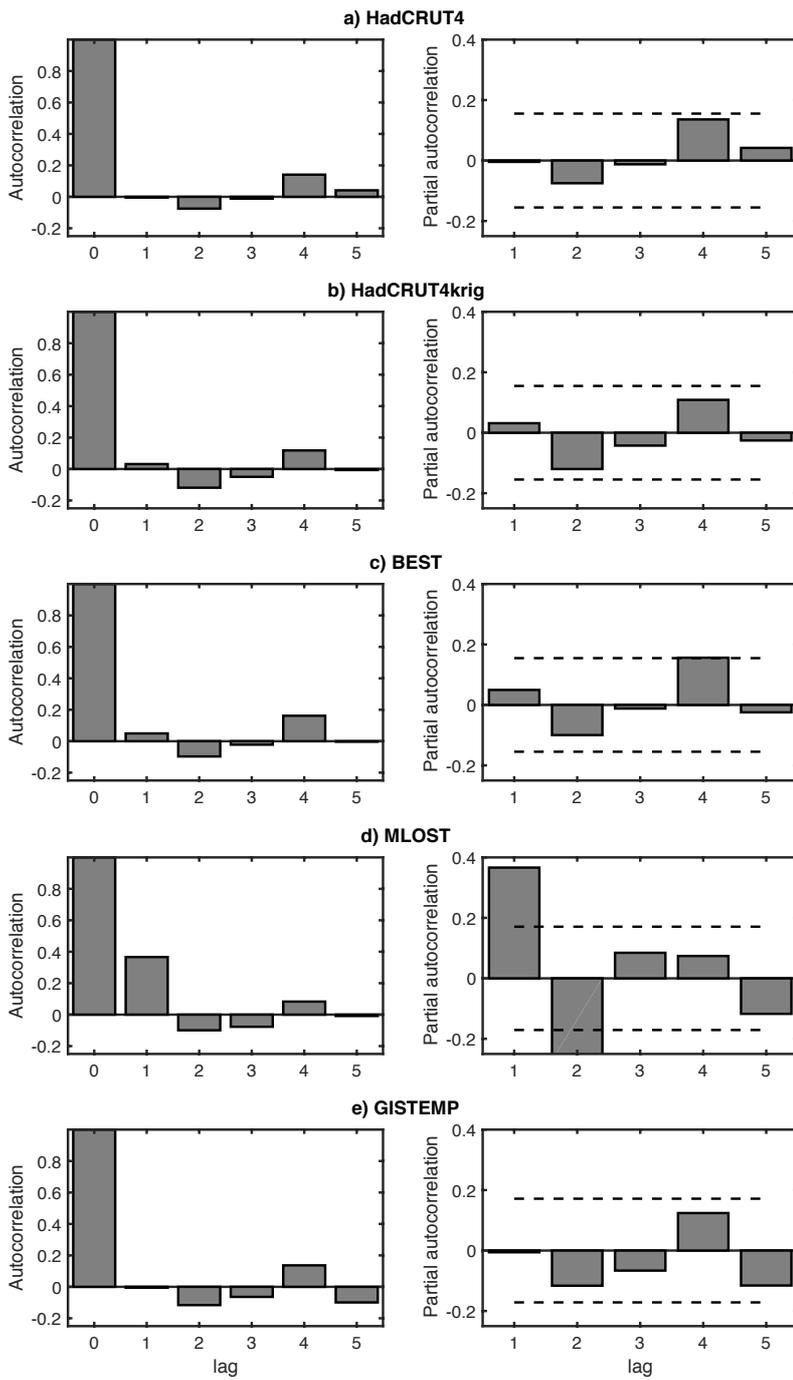
1149

1150 **Figure A3:** Number of change-points detected with EnvCpt, and STARS and BMCpt with
 1151 pre-whitening for each simulated scenario across 1,000 replications a) a linear trend, b) a
 1152 linear trend with first-order autocorrelation, c) a trend with three change-points in the
 1153 regression parameters, d) a trend with a change-point in the regression parameters and first-
 1154 order autocorrelation, e) a constant mean, f) a constant mean with first-order
 1155 autocorrelation, g) two change-points in the mean and h) two change-points in the mean with first-order
 1156 autocorrelation. The pre-whitening is performed using the using the MP and INV approaches
 1157 with a subsample size of 20.



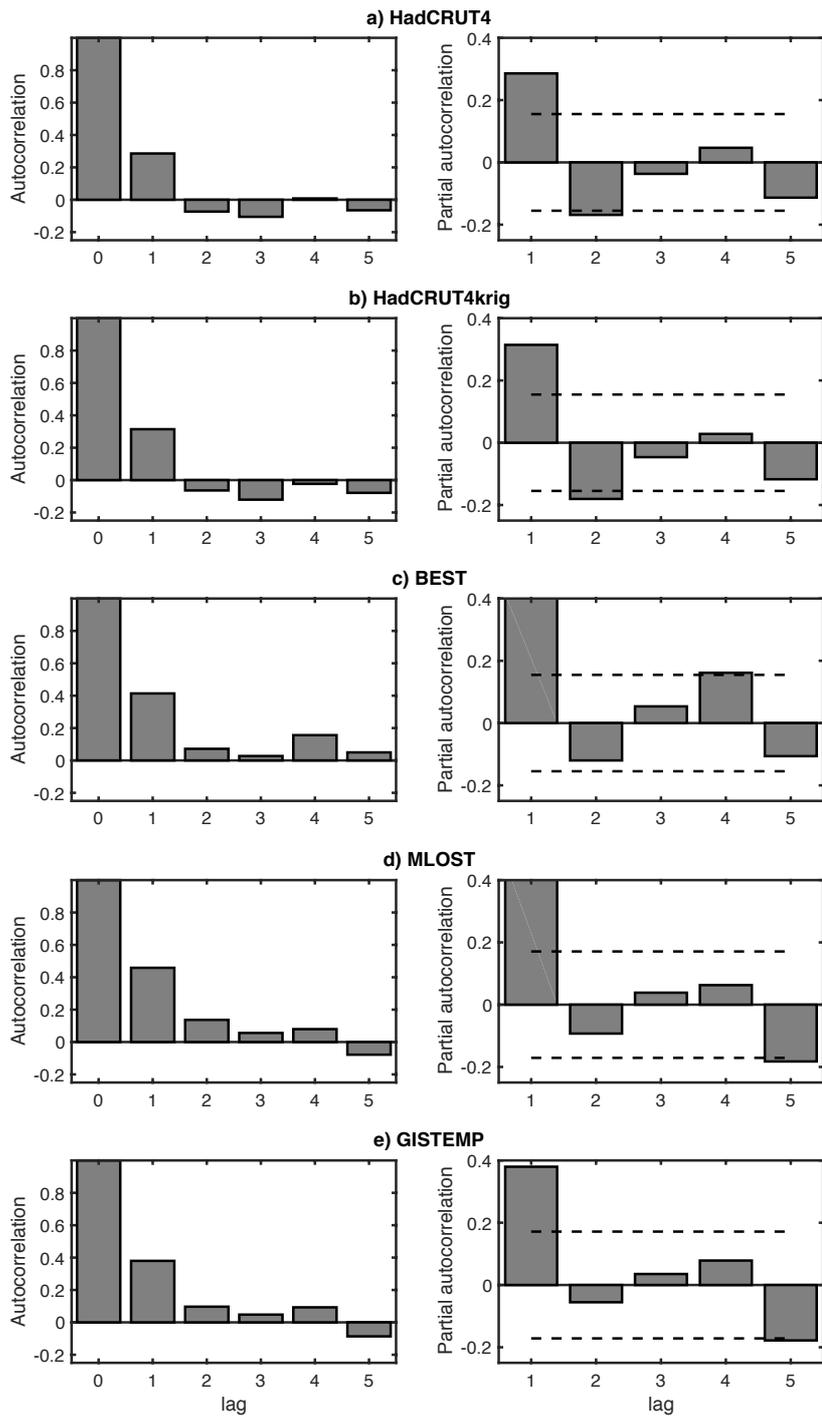
1158

1159 **Figure A4:** Density of change-point timings detected using EnvCpt, STARS and BMCpt
 1160 with pre-whitening for the two simulated scenarios with change-points and AR(1) across
 1161 1,000 replications a) a trend with a change-point in the regression parameters and first-order
 1162 autocorrelation and b) two change-points in the mean with first-order autocorrelation. The
 1163 pre-whitening is performed using the using the MP and INV approaches with a subsample
 1164 size of 20.



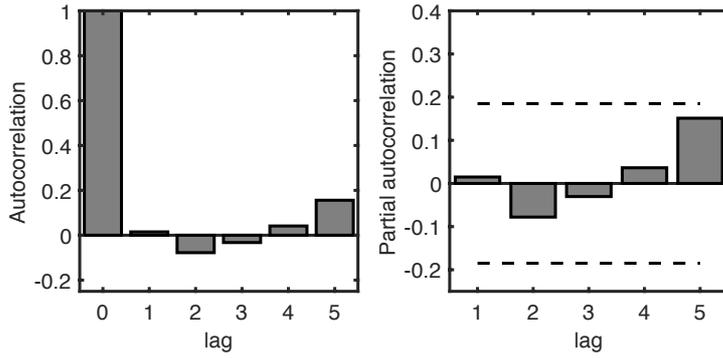
1165

1166 **Figure A5:** Autocorrelation and partial autocorrelation function of the residuals from the
 1167 Trend cpt + AR(1) model fitted to the global mean surface temperature datasets a)
 1168 HadCRUT4, b) HadCRUT4krig, c) BEST, d) MLOST and e) GISTEMP. Dashed lines
 1169 represent the 95% confidence intervals on the partial autocorrelation.



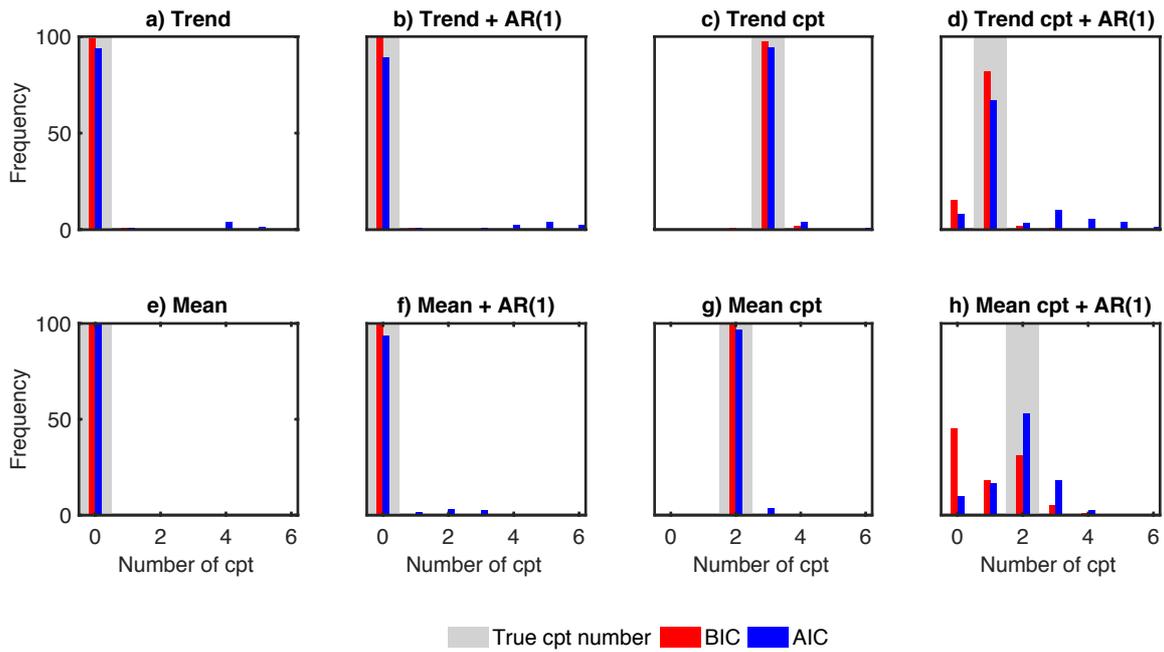
1170

1171 **Figure A6:** Autocorrelation and partial autocorrelation function of the residuals from the
 1172 Trend cpt model fitted to the global mean surface temperature datasets a) HadCRUT4, b)
 1173 HadCRUT4krig, c) BEST, d) MLOST and e) GISTEMP. Dashed lines represent the 95%
 1174 confidence intervals on the partial autocorrelation.



1175

1176 **Figure A7:** Autocorrelation and partial autocorrelation function of the residuals from the
 1177 Mean + AR(1) model fitted to the PDO. Dashed lines represent the 95% confidence intervals
 1178 on the partial autocorrelation.



1179

1180 **Figure A8:** Number of change-points detected with EnvCpt with either the Akaike
 1181 Information Criterion (AIC) vs the Bayesian Information Criterion (BIC) for each simulated
 1182 scenario across 1,000 replications a) a linear trend, b) a linear trend with first-order
 1183 autocorrelation, c) a trend with three change-points in the regression parameters, d) a trend
 1184 with a change-point in the regression parameters and first-order autocorrelation, e) a constant
 1185 mean, f) a constant mean with first-order autocorrelation, g) two change-points in the mean
 1186 and h) two change-points in the mean with first-order autocorrelation.