

1 **Continental-scale temperature variability in PMIP3 simulations and PAGES 2k regional**
2 **temperature reconstructions over the past millennium**

3
4 PAGES2k-PMIP3 group: Oliver Bothe¹, Michael Evans², Laura Fernández Donado³, Elena
5 Garcia Bustamante⁴, Joelle Gergis⁵, J. Fidel Gonzalez-Rouco³, Hugues Goosse^{6,a}, Gabi
6 Hegerl⁷, Alistair Hind^{8,9}, Johann Jungclaus¹, Darrell Kaufman¹⁰, Flavio Lehner^{11,12}, Nicholas
7 McKay¹⁰, Anders Moberg⁸, Christoph C. Raible¹², Andrew Schurer⁷, Feng Shi¹³, Jason E.
8 Smerdon¹⁴, Lucien von Gunten¹⁵, Sebastian Wagner¹⁶, Elliott Warren¹⁷, Martin Widmann¹⁷,
9 Pascal Yiou¹⁸ Eduardo Zorita¹⁶.

10 1. Max Planck Institute for Meteorology, Hamburg, Germany

11 2. Department of Geology & ESSIC, University of Maryland, USA

12 3. Institute of Geoscience (UCM-CSIC), Dpt. Astrophysics and Atmospheric Sciences, University
13 Complutense Madrid, Spain

14 4. Department of Physics. University of Murcia, Spain.

15 5. School of Earth Sciences, University of Melbourne, Australia

16 6. ELIC/TECLIM, Université Catholique de Louvain, Belgium

17 7. School of GeoSciences, University of Edinburgh, Edinburgh, United Kingdom

18 8. Department of Physical Geography, Stockholm University, Sweden

19 9. Department of Mathematics, Stockholm University, Sweden

20 10. School of Earth Sciences & Environmental Sustainability

21 Northern Arizona University

22 11. National Center for Atmospheric Research, Boulder, USA

23 12. Climate and Environmental Physics and Oeschger Centre for Climate Change Research,
24 University of Bern, Switzerland

25 13 Key Laboratory of Cenozoic Geology and Environment, Institute of Geology and Geophysics,
26 Chinese Academy of Sciences, 100029 Βειρωτινγ, Χητινα.

27 14. Lamont-Doherty Earth Observatory of Columbia University, Palisades, New York, USA

28 15. PAGES International Project Office, Falkenplatz 16, 3012 Bern, Switzerland

29 16. Institut for Coastal Research, Helmholtz Zentrum Geesthacht, Germany

30 17. School of Geography, Earth and Environmental Sciences, University of Birmingham, United
31 Kingdom

32 18. Laboratoire des Sciences du Climat et de l'Environnement, UMR 8212 CEA-CNRS-UVSQ, CE
33 Saclay l'Orme des Merisiers, 91191 Gif-sur-Yvette, France

34
35 ^a corresponding author

36 **Final Version: 27-11-2015**

37

38 **Abstract**

39 Estimated external radiative forcings, model results and proxy-based climate reconstructions
40 have been used over the past several decades to improve our understanding of the
41 mechanisms underlying observed climate variability and change over the past millennium.
42 Here, the recent set of temperature reconstructions at the continental-scale generated by the
43 PAGES 2k project and a collection of state-of-the-art model simulations driven by realistic
44 external forcings are jointly analysed. The first aim is to estimate the consistency between
45 model results and reconstructions for each continental-scale region over the time and
46 frequency domains. Secondly, the links between regions are investigated to determine whether
47 reconstructed global-scale covariability patterns are similar to those identified in model
48 simulations. The third aim is to assess the role of external forcings in the observed
49 temperature variations. From a large set of analyses, we conclude that models are in relatively
50 good agreement with temperature reconstructions for Northern Hemisphere regions,
51 particularly in the Arctic. This is likely due to the relatively large amplitude of the externally
52 forced response across northern and high latitude regions, which results in a clearly detectable
53 signature in both reconstructions and simulations. Conversely, models disagree strongly with
54 the reconstructions in the Southern Hemisphere. Furthermore, the simulations are more
55 regionally coherent than the reconstructions perhaps due to an underestimation of the
56 magnitude of internal variability in models or to an overestimation of the response to the
57 external forcing in the Southern Hemisphere. Part of the disagreement might also reflect large
58 uncertainties in the reconstructions, specifically in some Southern Hemisphere regions, which
59 are based on fewer paleoclimate records than in the Northern Hemisphere.

60

61 **1. Introduction**

62 The past millennium is an important period for testing our understanding of the
63 mechanisms that give rise to climate system variability (e.g., Masson-Delmotte et al., 2013).
64 Constraints on, and uncertainties in, external radiative forcings that drive climate change have
65 been extensively documented (e.g., Schmidt et al., 2011, 2012). Such radiative forcing data
66 sets can be used to drive climate simulations using the same model versions that are applied to
67 simulate future climate changes. This allows an evaluation of the relative importance of the
68 various forcings over time, while comparisons of past and future climate simulations place
69 20th-century climate variability within a longer context (e.g., Schmidt et al., 2014; Cook et al.
70 2015). Additionally, the availability of high quality paleoclimatic observations for the last
71 1000 years permits the reconstruction of regional, hemispheric and global scale climate
72 variability (e.g., Mann et al., 1999; Cook et al., 1999, 2004, 2010; Mann et al., 2009; Jones et
73 al., 2009; PAGES 2k Consortium, 2013, 2015; Masson-Delmotte et al., 2013; Neukom et al.,
74 2014). As a result, the past millennium has become a useful test case for evaluating climate
75 and earth system models used within the Intergovernmental Panel on Climate Change (IPCC)
76 fifth assessment report (Flato et al., 2013; Bindoff et al., 2013).

77 Paleoclimate reconstructions provides opportunities to test the fidelity of modelled
78 processes and their role in explaining past climatic variations. Reconstructions and
79 simulations can also be used jointly to evaluate estimates of climate sensitivity to external
80 radiative forcing (e.g., Hegerl et al., 2006; Braconnot et al., 2012; Masson-Delmotte et al.,
81 2013). Comparisons across many realizations of simulated climate are used to assess the

82 extent to which characteristic climate statistics are accurately simulated, as well as to
83 disentangle unforced and forced patterns (e.g., Hargreaves et al., 2013; Bothe et al., 2013ab;
84 Neukom et al., 2014; Coats et al., 2015a,b). Estimates of the unforced variability of the
85 climate system may be made from unforced simulations, or from the residual obtained when
86 the forced signal is removed from climate reconstructions, using realistically forced model
87 experiments (Schurer et al., 2013).

88 Furthermore, simulations can provide the basis for the design of observing network
89 arrays (Comboul et al., 2015). Simulation results also provide a test bed for paleoclimatic
90 reconstruction algorithms within so-called pseudo-proxy experiments (e.g., Zorita et al., 2003;
91 Hegerl et al., 2007; Smerdon, 2012; Lehner et al., 2012; Tingley et al., 2012; Wang et al.,
92 2014; Smerdon et al., 2015). All of these purposes, which are also pursued within the
93 historical period and with comparison to direct climate observations (Bindoff et al., 2013;
94 Ding et al., 2014), are potentially extended by the longer time interval made possible by
95 analysis over the past millennium.

96 However, obtaining unequivocal conclusions from the comparison between
97 reconstructions and simulation results over the past millennium remains difficult due to
98 uncertainties in climate and forcing reconstructions, the simplified world represented by
99 climate models, and the relatively weak forced signal in the pre-industrial part of the past
100 millennium compared to internal climate variability (e.g., Moberg, 2013). Reconstructions
101 and simulations are two different representations of the behaviour of the actual climate
102 system, and this creates multiple uncertainties in the task of intercomparison. Simulations
103 have uncertain forcings (Schmidt et al., 2011, 2012), and models contain parameterized or
104 uncertain representation of the physics, chemistry, biology and interactions within the climate
105 system (Flato et al., 2013). Furthermore, computational constraints impose a limited spatial
106 resolution or a deliberate omission of some known processes in order to perform simulations
107 at global scale that cover several centuries (e.g., Goosse et al., 2005; Schurer et al., 2013,
108 Phipps et al., 2014)

109 The uncertainty in paleoclimatic reconstructions is not always well understood either
110 and estimating its magnitude is challenging. For regional- to large-scale temperature
111 reconstructions, uncertainty can be caused by random or systematic error in the proxy
112 measurement, inadequate understanding of the proxy system response to environmental
113 variation, differences in fields derived from instrumental records selected to calibrate the
114 records, changes in the spatiotemporal and data type availability across the observational
115 network, and reconstruction methods (e.g., Jones et al., 2009; Smerdon et al., 2010; Smerdon,
116 2012; Emile-Geay et al., 2013; Evans et al., 2013; Wang et al., 2014; Comboul et al., 2015).

117 The non-climatic noise in reconstructions has a significant influence on model data
118 comparison. This may first have an impact on the variance of the reconstructed climatic signal
119 itself, although this is dependent on the actual choice of calibration method (e.g., Hegerl et al.,
120 2007; Christiansen et al., 2009; Mann et al., 2009; Smerdon et al., 2010; Smerdon, 2012).
121 Furthermore, the non-climatic noise can mask real relationships between climate variations in
122 different regions, or obscure the responses to forcing, which are clearer in models because of
123 the absence of this noise.

124 Acknowledging the considerable uncertainty in paleoclimatic reconstructions, the
125 earliest comparisons of past millennium simulations and reconstructions focused on
126 hemispheric- and global-scale changes, using a single, often simple, climate model driven by

127 globally uniform external radiative forcing estimates (e.g., Crowley, 2000; Bertrand et al.,
128 2002). Later, simulations with more comprehensive models (e.g., Gonzalez-Rouco et al.,
129 2006; Amman et al., 2007; Tett et al., 2007) refined the conclusions reached previously and
130 enabled regional and continental-scale analyses. They underscored the potential role of the
131 spatial distribution of some forcings, such as land use and of the dynamic response of the
132 atmospheric circulation (e.g., Luterbacher et al., 2004; Raible et al., 2006; Goosse et al., 2006;
133 Hegerl et al., 2011). Changes in the latter may be driven by the forcings (e.g., Shindell et al.,
134 2001; Mann et al., 2009) or be a signature of internal variability of the climate system (e.g.,
135 Wunsch, 1999; Raible et al., 2005).

136 State-of-the-art climate models reasonably simulate properties of internal variability,
137 such as teleconnection patterns or the probability of a particular event (e.g., Flato et al., 2013).
138 However, they are not expected to reproduce the part of the observed time trajectory that is
139 not directly constrained by external forcing because of the non-linear, chaotic nature of the
140 system (Lorenz, 1963). This makes model-data comparison a complex issue when using a
141 single simulation, because differences between model results and reconstructions may be due
142 to a model or reconstruction bias, but may also simply reflect a different sample of internal
143 variability (defined here as the fraction of climatic variability that is not due to changes in
144 external forcings).

145 Indeed, comprehensive climate models have their own internal climate variability and,
146 if a model represents the real world in a satisfactory way, the observed trajectory would just
147 be one among all potential model realizations. The issue may be addressed by analysing an
148 ensemble of simulations, which provides information on the range that can be simulated by
149 one single model (e.g., Goosse et al., 2005; Yoshimori et al., 2005; Jungclaus et al., 2010;
150 Moberg et al., 2015) or a set of models (e.g., Jansen et al., 2007; Lehner et al., 2012;
151 Fernandez-Donado et al., 2013, Bothe et al., 2013b). The reconstruction has then to be
152 compatible with this range, at least when considering all the uncertainties, to claim
153 consistency between simulations and reconstructions, whereby such a compatibility can be
154 defined in various ways, as discussed below.

155 Fernandez-Donado et al. (2013) reviewed results from 26 climate simulations with eight
156 atmosphere-ocean general circulation models (AOGCMs), reflecting the state of modelling
157 before the CMIP5/PMIP3 (Coupled Model Intercomparison Project Phase 5/Paleoclimate
158 Model Intercomparison Phase 3). These pre-CMIP5/PMIP3 simulations were driven by a
159 relatively wide range of choices for boundary conditions and forcing agents. For the Northern
160 Hemisphere surface temperature variations, Fernandez-Donado et al. (2013) found an overall
161 agreement within the temporal evolution, but still noted discrepancies between simulations
162 and hemispheric and global temperature reconstructions. For example, the period between
163 around 850 and 1250 CE is warmer in the reconstructions than in the simulations (see also
164 Jungclaus et al., 2010; Goosse et al., 2012b; Shi et al., 2013).

165 Additionally, a comparison of the simulated changes in the temperature fields from this
166 warm period and the colder period around 1450-1850 showed little resemblance to the field-
167 reconstruction by Mann et al. (2009), but the spatial reconstructions themselves have
168 significant uncertainties (e.g., Wang et al., 2015). These two relatively warm and cold periods
169 are often referred to as the Medieval Climate Anomaly (MCA), and the Little Ice Age (LIA),
170 respectively, although their exact timing has been debated and the adequacy of their names
171 has been questioned (e.g., Jones and Mann, 2004; PAGES 2k Consortium, 2013).

172 The assessment of information from paleoclimate archives (Masson-Delmotte et al.,
173 2013) in the IPCC fifth assessment report partly followed the approach applied by Fernandez-
174 Donado et al. (2013). Masson-Delmotte et al. (2013) included a preliminary analysis of the
175 more recent CMIP5/PMIP3 “past1000” simulations, which were coordinated more closely
176 than previous experiments, particularly in regard to the choices of forcings (Schmidt et al.,
177 2011, 2012). They came to similar conclusions as Fernandez-Donado et al. (2013): the
178 reconstructed MCA warming is greater than simulated, but not inconsistent within the large
179 uncertainties.

180 Agreement between paleoclimate reconstructions and simulations has also been
181 assessed by compositing the response to individual forcing events (e.g., Hegerl et al., 2003;
182 2011; Luterbacher et al., 2004; Stenchikov et al., 2006; Masson-Delmotte et al., 2013). The
183 reconstructed and simulated response to volcanic forcing agrees in magnitude on multi-
184 decadal time scales. Detailed comparisons of observations around the 1815 Tambora eruption
185 indicate that the simulated cooling is larger than in instrumental observations or in
186 reconstructions (Brohan et al., 2012), but a significant part of the discrepancy might be due to
187 forcing uncertainties.

188 For the solar forcing, direct comparisons between simulations and reconstructions are
189 inconclusive regarding whether simulations that use either moderate or weak variations of
190 total solar irradiance provide generally better agreement with reconstructions (Masson-
191 Delmotte et al.; 2013; Fernandez-Donado et al., 2013). This has been confirmed at
192 hemispheric and regional-scales by Hind and Moberg (2013) and Moberg et al. (2015), using
193 appropriately designed statistical tests of temporal correlation and quadratic distance between
194 reconstructions and simulations (Sundberg et al., 2012).

195 The cause of past climate change in the Northern Hemisphere, specifically the
196 contribution by individual forcings to a climatic event, can be estimated using detection and
197 attribution techniques. These techniques allow for the possibility that the reconstructions
198 contain forced signals of larger or smaller magnitude than simulated (e.g., due to forcing
199 uncertainty, uncertainty in a models transient response, or uncertainty in calibration of
200 reconstructions). The results show that the response to volcanic eruptions can be clearly
201 detected in reconstructions, consistent with epoch analysis results, and also confirm that the
202 signal is generally larger in magnitude in the simulations (Hegerl et al., 2003; 2007; Schurer
203 et al., 2013), although the discrepancy may be within the range of volcanic forcing
204 uncertainty. The response to solar forcing cannot be reliably separated from internal
205 variability, but very high solar forcing such as that reconstructed by Shapiro et al. (2011)
206 needs to be significantly scaled down to match reconstructions even given large
207 reconstruction uncertainties (Schurer et al., 2014). Within the LIA, detection and attribution
208 methods show that volcanic forcing is critical for explaining the anomalous cold conditions
209 (Hegerl et al., 2007; Miller et al., 2012; Lehner et al., 2013; McGregor et al., 2015) and that
210 there is also weak evidence for a contribution from a small but long-lived decrease in CO₂
211 concentration (e.g., MacFarling-Meure et al., 2006; Schurer et al., 2014).

212 The studies mentioned above mainly focused on the Northern Hemisphere, because a
213 larger number of paleoclimatic observations and reconstructions are available there. However,
214 several recent studies assessed differences in inter-hemispheric connections (Goosse et al.,
215 2004; Neukom et al., 2014), Southern Hemisphere climate variability (Phipps et al., 2014),
216 regional temperature variability (Luterbacher et al., 2004; Hegerl et al., 2011; Goosse et al.,

217 2012a; Gergis et al., 2015; Shi et al., 2015) and Southern Hemisphere circulation features
218 (Wilmes et al., 2012; Abram et al., 2014; Tierney et al., 2015).

219 In particular, the recent consolidation of Southern Hemisphere paleoclimate data
220 (Neukom and Gergis, 2012) led to the comparison of a hemispheric temperature comparison
221 with a suite of 24 climate model simulations spanning the past millennium (Neukom et al.,
222 2014). This study reported considerable differences in the 1000-year temperature
223 reconstruction ensembles from the Northern and Southern Hemispheres. An extended cold
224 period (1590s–1670s CE) was observed in both hemispheres, while the current (post-1974)
225 warm phase is found to be the only period of the past millennium where both hemispheres
226 experienced simultaneous warm anomalies (Neukom et al., 2014). Their analyses also
227 suggested that the simulations underestimate the influence of internal variability in the ocean-
228 dominated Southern Hemisphere (Neukom et al., 2014).

229 While several studies have provided valuable advances in our understanding of
230 hemispheric-scale climate dynamics, this brief overview indicates that observed and simulated
231 paleoclimate variations at regional and continental scales have not been thoroughly compared
232 up to now. This was the goal of a workshop joining the PAGES 2k and PMIP3 communities in
233 Madrid (Spain) in November 2013, using a recent set of continental-scale temperature
234 reconstructions (PAGES 2k Consortium, 2013) and a collection of state-of-the-art model
235 simulations driven by realistic external forcings (Schmidt et al., 2011, 2012). On the basis of
236 the discussions held during this workshop, the aim of this study is to systematically estimate
237 the consistency between the simulated and reconstructed temperature variations at the
238 continental scale and evaluate the origin of observed and simulated variations. This study is
239 motivated by the following key science questions:

240 1/ Are the statistical properties of surface temperature data for each individual
241 continent-scale region consistent between simulations and reconstructions?

242 2/ Are the cross regional relations of temperature variations similar in reconstructions
243 and models?

244 3/ Can the signal of the response to external forcing be detected on continental scale
245 and, if so, how large are these signals?

246 Section 2 first presents a brief overview of the PAGES 2k reconstructions and
247 simulations analysed here. In addition to a selection of PMIP3 simulations, some numerical
248 experiments that did not follow the PMIP3 protocol were also analysed, mainly to include
249 model runs with larger solar forcing amplitude. We use several statistical methods to achieve
250 robust results in answering the key science questions above. They are listed at the end of
251 section 2. Each methodology is briefly described when it is applied while some specific
252 implementation information is provided in supplement sect. S2. In section 3, each continental-
253 scale region is studied separately to determine whether the reconstructed and simulated time
254 series have similar characteristics, in terms of the magnitude and timing of the observed
255 changes as well as the spectral distribution of the variance. Section 4 investigates whether the
256 inter-regional patterns of temperature variability are similar in the reconstructions and
257 simulations. The role of the external forcings in producing the observed variations is
258 presented in section 5. Section 6 provides a discussion of our results, their limitations and how
259 our conclusions compare to previous studies. Finally, section 7 summarizes the main findings

260 and provides perspectives for future developments. Several additional analyses are provided
261 as supplementary material for completeness and further reference.

262

263 **2. Data and methods**

264 **2.1 PAGES 2k reconstructions**

265 The PAGES 2k Consortium (2013) generated temperature reconstructions for seven
266 continental-scale regions (Fig. 1). The proxy climate records found to be best suited for
267 reconstructing annual or warm-season temperature variability within each continental-scale
268 region were identified. Expert criteria for the adequacy of proxies were *a priori* specified
269 (PAGES 2k Consortium, 2013). The resulting PAGES 2k dataset includes 511 time series
270 from different archives including tree rings, pollen, corals, lake and marine sediment, glacier
271 ice, speleothems, and historical documents. These data record changes in biological or
272 physical processes and are used to reconstruct temperature variations (all data are archived at:
273 https://www.ncdc.noaa.gov/cdo/f?p=519:2:0:::P1_study_id:12621).

274 The PAGES 2k reconstructions have annual resolution in all regions except North
275 America, which has one 780-year-long tree-ring-based reconstruction (back to 1200 CE with
276 10-year resolution) and one 1400-year-long pollen-based reconstruction (back to 480 CE with
277 30-year resolution). These latter two reconstructions therefore are smoothed differently and
278 they are either excluded from the analysis or treated in slightly different ways in some
279 comparisons. The reconstruction for the Arctic region used in this study is based on a revised
280 version (v 1.1) of the PAGES 2k dataset (McKay and Kaufman, 2014).

281 Each regional group tailored its own procedures to their local proxy records and
282 regional calibration targets (PAGES 2k Consortium, 2013). Thus, each continental-scale
283 temperature reconstruction was derived using different statistical methods. In short, most
284 groups used either a scaling approach to adjust the mean and variance of a predictor
285 composite to an instrumental target, or a regression-based technique to extract a common
286 signal from the predictors using principal components or distance weighting. Thus, some of
287 the observed region-to-region differences between simulations and reconstructions might be
288 due to the differences in reconstruction methods. Nevertheless, alternative reconstructions for
289 all regions based on exactly the same statistical procedures were also produced and were
290 found to be similar to the PAGES 2k temperature reconstructions provided by each group
291 (PAGES 2k Consortium, 2013). Each regional group also used individually selected
292 approaches to assess the uncertainty of their temperature reconstructions, designed to quantify
293 different aspects of the uncertainty. For example, some regions primarily quantified
294 uncertainties associated with the set of records used in the reconstruction and their agreement
295 through time, which can reflect within-region variability as well as uncertainty (Arctic, North
296 American tree rings). Other regions focused on uncertainties associated with how closely the
297 proxy resembles temperatures (Asia, Antarctica, Europe, North American pollen), and some
298 regions incorporated both of these types of uncertainties (Australasia, South America).
299 Uncertainty estimates in all of the regions except for Antarctica vary through time depending
300 on the set of records available for any given interval and their agreement. All uncertainty
301 estimates that assess how well the proxy data reproduce observed temperatures are based on
302 the assumption that the modern proxy-temperature relation is stationary into the past, and that

303 the agreement between proxy data and temperature on short timescales can be used to infer
304 uncertainty at lower frequencies.

305 **2.2 Climate model simulations**

306 The climate model simulations used in this study are listed in Table 1, summarizing
307 model specifications such as resolution, forcing applied to the transient simulations, and
308 length of pre-industrial control simulations (piControl). These simulations include
309 contributions to the third Paleoclimate and the fifth Coupled Modelling Intercomparison
310 Projects (PMIP3, Braconnot et al., 2012; CMIP5, Taylor et al., 2012) from six models
311 (CCSM4, CSIRO-Mk3L-1-2, GISS-E2-R, HadCM3, IPSL-CM5A-LR, MPI-ESM-P), as well
312 as a more recent simulation with CESM1, and the COSMOS pre-PMIP3 ensemble with
313 ECHAM5/MPIOM (see also Table S1).

314 The experiments were selected among available pre-PMIP3 and PMIP3 simulations on
315 the basis of specific criteria: The conditions were (i) they run continuously from 850-2000
316 CE; (ii) they include at least solar, volcanic aerosol, and greenhouse gas forcing (S, V, G in
317 Table 1); (iii) they use a plausible solar forcing reconstruction with an amplitude within the
318 range that is consistent with recent understanding (iv) they do not display a large unphysical
319 drift over the simulated period.

320 PMIP3 simulations all comply with criteria (ii) and (iii) as they follow the
321 recommendation of Schmidt et al. (2011) by using an increase in total solar irradiance (TSI)
322 from the late Maunder Minimum period to the present day of $\sim 0.10\%$. Nevertheless, some
323 PMIP3 simulations were excluded from the analysis, as the simulations presented clear
324 incompatibilities with the rest of the ensemble. For instance, the MIROC simulation displays
325 a trend in the global annual mean temperature over the whole millennium that is not
326 compatible with the present understanding of the past millennium climate. It has been
327 considered here as a likely model artefact that could also affect regional and seasonal
328 temperatures in unknown ways. Contrary to the GISS model, this drift is not clearly
329 understood and no control run is available to statistically correct it. The simulation with bcc-
330 csm-1 was discarded because of potentially unphysical large anomalies in some regions.
331 FGOALS-g1 was not use due to the unavailability of a continuous run from 850 to 2000, as
332 the so-called 'past1000' simulation covers only the years 850-1850 under the PMIP protocol.

333 Most non-PMIP3 simulations did not comply with at least one the criteria above.
334 Nevertheless, experiments performed with two models (ECHAM5/MPIOM and CESM1)
335 follow all of them. They include simulations with a stronger solar forcing than in the PMIP3
336 ensemble. A three-member ensemble with ECHAM5/MPIOM uses a TSI reconstruction with
337 an increase of $\sim 0.24\%$ (COSMOS E2), while CESM1 uses a TSI reconstruction with an
338 increase of $\sim 0.20\%$. No simulation used in this study incorporates the much larger increase of
339 $\sim 0.44\%$, suggested by Shapiro et al. (2011), which results in simulations that are inconsistent
340 with reconstructed large-scale temperatures (Feulner, 2011; Schurer et al., 2014). The
341 COSMOS simulations deviate from the PMIP3 protocol because they included an interactive
342 carbon cycle with CO_2 concentration as prognostic variable. While simulated and
343 reconstructed CO_2 evolution diverge during some periods, the differences have only a
344 marginal effect on simulated temperatures (Jungclaus et al., 2010).

345 Consequently, the group of simulations analysed here is not strictly based on the PMIP3
346 ensemble. Nevertheless, as we use a majority of PMIP3 simulations and additional

347 simulations that follow an experimental design similar to PMIP3, we will keep the reference
348 to PMIP3 for simplicity.

349 The variable extracted from the simulation outputs is the monthly mean surface air
350 temperature (labelled ‘tas’ in the Climate Model Output Rewriter framework of CMIP5).
351 These temperature fields were then used to create area-averaged time series that matched the
352 domain and seasonal window of each of the PAGES 2k regional reconstructions (see
353 Supplement Sect. S1).

354 **2.3 Statistical methods**

355 Several climate model–paleoclimate data comparison and analysis methods are used in
356 this study to verify the robustness of the results generated by each method and to provide a
357 comprehensive guide for future work. Model-data comparisons need to account for
358 uncertainties in climate reconstructions, in forcing reconstructions and in the response to
359 forcings in model simulations. These approaches also must recognize that the real climate,
360 and hence the reconstructions, and individual climate model simulations include their own
361 individual realizations of internally generated variability. Therefore, perfect agreement
362 between model simulations and data can never be expected when directly comparing time
363 series.

364 The first group of methods is focused on the first question raised in the introduction.
365 The goal is to assess if temperature reconstructions have similar statistical properties
366 compared to simulations. This is initially done by simple analysis of the time series, such as
367 estimates of the variance (section 3.1). The spectral properties are then analysed (section 3.2)
368 before the probabilistic and climatological consistency (section 3.3) and the skill of the
369 various simulations (section 3.4). The second question dealing with the cross regional
370 variations of temperatures is addressed by discussing the correlation between regions
371 (sections 4.1 and 4.3) and through a principal component analysis (section 4.2). Finally, the
372 third question about the role of the forcing is studied by means of a superposed epoch analysis
373 (section 5.1), by applying a statistical framework involving correlation and distance metrics
374 (section 5.2) and detection and attribution techniques (section 5.3). For more details on those
375 methods, see supplement sect. S2.

376 In the majority of the analyses presented in this manuscript, anomalies compared to
377 the mean over the whole period covered are used and the time series are smoothed or
378 temporally averaged, using either a 23-point Hamming filter or non-overlapping 15-year
379 averages, depending on the requirements of the various techniques (both methods give a
380 similar degree of low-pass filtering). This is motivated by the relatively weaker skill of some
381 reconstructions to replicate observed records on interannual time scales (Cook et al., 2004;
382 Esper et al., 2005; D’Arrigo et al., 2006) and by the fact that the main focus here on decadal
383 to centennial timescales. The full period analysed is 850-2005 CE, although different periods
384 are chosen for some analyses because of data availability, the choice of the temporal filtering,
385 other technical restrictions, or to analyze sub-periods.

386

387 **3. Regional analysis**

388 To begin, the agreement between simulations and reconstructions for individual
389 regions is described qualitatively, using a simple visual comparison of the time series, and

390 then quantitatively by calculating spectra, consistency and skill metrics. The correlations
391 between the time series are presented in the supplementary material (Figure S1 and
392 supplementary text S3). Overall, the analyses in this section illustrate the potential of
393 identifying common signals in both data sets. The different diagnostics are presented here
394 separately whereas the conclusions derived from the results of the different analyses are
395 compared and discussed in more detail in section 6.

396 **3.1 Observed and simulated time series**

397 Figure 1 shows the regional time series in the forced simulations with each regional
398 temperature reconstruction. To the right of each time series graph, the magnitude of variability
399 of unforced simulated temperatures is illustrated by calculating the standard deviation of pre-
400 industrial control simulations in each model. The unforced variability is generally similar in
401 all models in all the regions, with weaker amplitudes in Australasia and Asia. Note that some
402 regions cover only land areas while others have an oceanic fraction (see Supplement Sect.
403 S1), with a potential impact on the magnitude of the estimated variability.

404 Most reconstructions show a tendency of a gradual cooling over the millennium,
405 followed by recent warming. Notable common features among regions on decadal timescales
406 are the pronounced negative anomalies related to large tropical volcanic eruptions in the
407 simulations. This is most obvious for the eruptions in the 1250s, 1450s and 1810s. Among the
408 temperature time series, a larger response to volcanic eruptions is noticeable in the CESM,
409 MPI and CCSM4 simulations. The regional temperature reconstructions rarely capture the
410 first two of these anomalies or only register them at smaller amplitudes. Only the early 19th
411 century eruptions are clearly reflected in many regions, and are most pronounced in the
412 Northern Hemisphere reconstructions. The reconstruction for Europe also shows a negative
413 anomaly coinciding with the effect of the 1450s eruption, with an amplitude comparable to
414 that seen in some of the simulations.

415 Figure 1 suggests that the temperature reconstructions show slightly more centennial to
416 multi-centennial variability than the models over the full period with stronger long-term
417 trends, while several model results indicate a stronger recent warming compared to some of
418 the reconstructions. The reconstruction uncertainty bands provided with the original PAGES
419 2k reconstructions encompass the simulated series with few exceptions, in particular the
420 Arctic and North America during the 1250s. The published uncertainty estimates have been
421 calculated using different methods for the various continental-scale regions, as detailed in the
422 supplementary material of PAGES 2k Consortium (2013). Furthermore, those uncertainties
423 are only valid at the original temporal resolution, which is annual in all cases except for North
424 America. It is expected that the reconstruction uncertainty decreases at lower resolution, or
425 after smoothing as in our case. This is consistent with the lower uncertainty ranges for the
426 low-resolution pollen-based reconstruction.

427 However, estimating the reduction of the uncertainty due to smoothing is not
428 straightforward (e.g., Moberg and Brattström, 2011; Franke et al., 2013) as the resulting
429 uncertainty magnitude is also dependent on autocorrelation of the non-climatic noise in proxy
430 data. The extreme hypothesis, considering that the error is constant in time and that the errors
431 are uncorrelated would lead to a decrease proportionally to 1 over the square root of the
432 number of samples included in the average. For a smoothing similar to 15-year averaging, as
433 performed herein, the approximation that likely leads to an underestimation of the
434 uncertainties would correspond to a decrease by a factor of about 4 compared to the original

435 error estimate. This suggests very small errors for most reconstructions. In this case, the major
436 discrepancies between the reconstructions and model results would occur at the same time as
437 mentioned above; however, periods when the models are out of the range of the
438 reconstruction uncertainty bands would be more common at the decadal scale.

439 For North America, the long term multi-centennial trend appears to be similar between
440 the pollen based reconstruction and simulations, except for the last ~200 years, when some
441 simulations show much stronger warming than is present in the reconstruction. This warming
442 feature is somewhat stronger in the tree-ring based reconstruction than in the pollen-based
443 reconstruction, but is nevertheless weaker than in some simulations. The COSMOS
444 simulations appear to be collectively colder than this reconstruction in the late 20th century.
445 Although the European temperature reconstruction and simulated series disagree substantially
446 in some parts of the 12th century and for the last ~200 years, there are otherwise strong
447 similarities, particularly during periods of large volcanic eruptions. Simulated and
448 reconstructed Arctic series show large decadal to centennial variability, but the timing of these
449 variations do not agree well. Therefore, simulations are often outside the reconstruction's
450 uncertainty range. Consistently, there is a large multi-model ensemble-spread but also single-
451 model ensemble spread as illustrated by the COSMOS simulations. CESM, CCSM4, and
452 IPSL show a strong recent warming and strong volcanic cooling.

453 Simulated and reconstructed temperatures show only weak long-term trends in Asia, but
454 decadal variability appears to be larger in the reconstruction. Simulations generally differ
455 from the reconstruction in the last 200 years and show either much weaker or much stronger
456 trends. In Australasia, the weak forced variability common to all simulations may be due to
457 the large spatial extent of the domain, which includes large oceanic areas that may dampen the
458 forced high-frequency variability. For the recent warming, the trends in CESM, CCSM4, IPSL
459 and the COSMOS simulations are considerably stronger than the Australasian temperature
460 reconstruction. The temperature reconstruction for South America is often near the upper or
461 lower limit of the simulation ensemble range and displays more centennial-scale variability
462 than the simulations. In Antarctica, the reconstruction has a clear long-term negative trend and
463 only a modest warming in the 20th century while the simulations show nearly no long-term
464 cooling but agree on the warming onset in the beginning of the 20th century.

465 **3.2 Spectral analysis**

466 Next, we consider the agreement between simulated and reconstructed temperature data
467 in terms of their spectral densities, which show how temperature variances are distributed
468 over frequency (Fig. 2, see also Fig. S2). Spectra were computed using the multi-taper method
469 (Thomson, 1982; Percival and Walden, 1993), with its so-called time-bandwidth product
470 being set to four. Consequently, each calculated spectrum is an average of seven statistically
471 independent spectrum estimates. Spectra for the reconstructions are illustrated with their 95%
472 confidence intervals, while model spectra are plotted with single lines. The analysis is made at
473 the original time resolution using all existing data points in the time window 850-2005.

474 The degree of agreement between model and reconstruction spectra differ substantially
475 between regions, with the Arctic showing the best agreement at all frequencies and South
476 America showing the worst. In the latter, most model spectra lie in the reconstruction
477 confidence interval only in a narrow frequency band corresponding to about 100- to 150-year
478 periods. The agreement is generally good for the Arctic, Europe and Asia and at multi-decadal
479 timescales (20-50 years) for many regions. Nevertheless, many models have systematically

480 less variance in the 50- to 100-year band and most models have more variance than the
481 reconstructions at higher frequencies.

482 Pronounced differences of high-frequency variance is seen for all Southern Hemisphere
483 regions. In particular, the pre-PMIP3 COSMOS simulations show significantly too much
484 variance at timescales of 3 to 5 years for Australasia and to a lesser degree for South America
485 and Antarctica. This property has previously been related in regions with strong influence
486 from tropical Pacific variability to this model's ENSO variability (Jungclauss et al., 2006;
487 Fernández-Donado et al., 2013). Most model spectra for North America lie within the
488 confidence interval of the tree-ring based reconstruction spectrum, although several models
489 have somewhat less variance than this reconstruction at periods longer than 50 years. The
490 North America pollen-based reconstruction behaves as a roughly 150-year low-pass filtered
491 series and has significantly less variance than the corresponding tree-ring-based record at all
492 frequencies for which both spectra are defined.

493

494 **3.3 Consistency estimate**

495 The probabilistic and climatological consistency of PMIP3 simulations and PAGES 2k
496 reconstructions was assessed following the framework of Annan and Hargreaves (2010; and
497 references therein; Hargreaves et al., 2011, 2013) and Marzban et al. (2011), respectively. The
498 current application is based on Bothe et al. (2013a,b). The underlying null hypothesis follows
499 the paradigm of a statistically indistinguishable ensemble (Annan and Hargreaves, 2010;
500 Rougier et al., 2013), i.e. the validation target, represented here by the temperature
501 reconstructions, and the model simulations are samples from a common distribution and are
502 therefore exchangeable.

503 Climatological consistency refers to the similarity of the climatological probability
504 distributions of reconstructions and of simulations over a selected period, either the whole
505 millennium or sliding sub-periods. We analyse climatological consistency by comparing
506 individual simulated series with the target (i.e., the reconstructions) to identify deviations in
507 climatological variance and possible biases between them. To achieve this goal, Marzban et
508 al. (2011) proposed the use of residual quantile-quantile (r-q-q) plots that should be
509 approximately flat for consistent series (supplement sect S2.1).

510 Probabilistic consistency refers to the position of the reconstruction in the range
511 spanned by the ensemble of simulations. Histograms of the ranks should be flat under
512 exchangeability (supplement sect S2.1), i.e. estimated frequencies of the verification target
513 and the ensemble agree if the simulation ensemble is probabilistically consistent with the
514 temperature reconstructions (Murphy, 1973).

515 As there are large uncertainties in paleoclimate reconstructions, it is necessary to take
516 into account these uncertainties in the evaluation of the consistency of the ensemble of
517 climate model simulations (Anderson, 1996). This is achieved by inflating the model
518 simulations results by adding noise with amplitudes that are proportional to published
519 uncertainty estimates from the original temperature reconstructions.

520 We assess probabilistic and climatological consistency based on non-overlapping 15-
521 year averages centred on the full period considered, except for the Northern American
522 temperature reconstruction where non-overlapping 30-year averages are used for the pollen-

523 based reconstruction, and 10-year averages for the tree-ring-based reconstruction. The results
524 are presented on Figure 3, Figure S3 and Figure S4 for all regions.

525 The regions selected for Figure 3 are chosen to provide a contrasting example. Two
526 estimates of the uncertainties are used. First, the uncertainties provided with the original
527 reconstruction are applied, which is an overestimation for the smoothed time series. Second,
528 at the other extreme, the uncertainties are assumed to be equal to zero and are thus known to
529 be underestimated. A third estimate of the uncertainty is provided in the supplementary
530 figures, using an uncertainty measure equal to the one provided in the original publication
531 divided by a factor $\sqrt{15}$ to account for the smoothing (see section 3.1). This leads to results
532 that are generally very similar to the case where uncertainty is assumed to be zero.

533 The simulations in most cases lack climatological consistency with the reconstructions
534 (Fig. 3 and Fig. S3). The simulated quantiles can deviate strongly from the reconstructed
535 quantiles. Specifically, the simulated distributions are generally over-dispersive when using
536 the original estimates of uncertainties. The differences are much smaller when uncertainties in
537 reconstructions are neglected, although extremes often remain overestimated. The Arctic and
538 the North American tree-ring based reconstruction are exceptions as some simulations are
539 climatologically consistent with the reconstruction and display only small differences between
540 simulated and reconstructed quantiles for all estimates of the uncertainty. Consistency is
541 reduced for those simulations that show larger variability (recall Figure 1) as is the case of the
542 CCSM4 and CESM models.

543 In agreement with the climatological assessment, the simulated results generally lack
544 probabilistic consistency with the reconstructions when the original uncertainty is considered
545 (Fig. 3 and Fig. S4). The target data are too often in the central ranks, indicating that the
546 probabilistic distribution of the ensemble is too wide and shows significantly over-dispersive
547 spread deviations. The only exception is the North American region using the tree-ring based
548 reconstruction. The most prominent differences are found in the Antarctic region where the
549 simulation ensemble spread deviates considerably from reconstructed temperatures (Fig. S4),
550 but strong ensemble spread deviations relative to the pollen reconstruction for North America
551 are also evident.

552 This assessment of the probabilistic consistency strongly depends on the estimate of the
553 uncertainty of the reconstruction. If we do not add noise to the model time series to reflect
554 error in reconstructions before the ranking and thereby neglect reconstruction uncertainty, or
555 if we assume a strong reduction of the error in reconstruction at the decadal time scale
556 because of the smoothing, the ensemble appears to be consistent with a number of regions or
557 even under-dispersive for others. However, ignoring the uncertainty in such a manner may
558 lead to an overconfident assessment of consistency between simulation ensemble and
559 reconstruction. Nevertheless, because the uncertainties are not well known, over-dispersion
560 does not necessarily weaken the reliability of the ensemble relative to the target, but instead
561 may highlight insufficiently constrained uncertainties in the reconstruction.

562 **3.4 Skill estimate**

563 The skill of the simulations is assessed using a metric introduced by Hargreaves et al.
564 (2013). The idea of skill stems from weather forecasting and refers to the ability of a
565 simulation to represent a target better than some simple reference values. For instance, in
566 weather forecasting, a standard reference is to assume no change compared to initial

567 conditions (i.e., persistence). A forecast has a positive skill if it is closer to the observed
 568 changes than this simple reference. The skill S , as in Hargreaves et al. (2013), is then:

$$569 \quad S = 1 - \sqrt{\frac{\sum (F_i - O_i)^2 - \sum e_i^2}{\sum (R_i - O_i)^2 - \sum e_i^2}} \quad (1)$$

570 where F_i is the simulation result at each data point, O_i is the reconstruction data, R_i is the
 571 reference (for instance a constant climate here) and e_i is uncertainty of the target. The square-
 572 root expression becomes undefined when either the actual simulation or the reference is better
 573 than the upper possible agreement level indicated by the errors. Uncertainty estimates are
 574 derived from the originally reported uncertainties in regional temperature reconstructions
 575 given by PAGES 2k Consortium (2013). If reconstructed error estimates are realistic, we do
 576 not expect the simulations to fit the target better than these uncertainty estimates. As for the
 577 consistency analyses, the skill analysis is calculated using temperature anomalies from the
 578 long-term averages within each analysis period.

579 Figure 4 presents the skill for the Arctic and Antarctica, as an example, with the other
 580 PAGES 2k regions displayed in Figure S5. In this estimate, we use a no-change reference
 581 forecast (i.e., the reference is the climatology) as there is no clear *a priori* evidence that the
 582 climate at one particular time during the past millennium is warmer or colder than the mean.
 583 Positive values suggest that the simulations is in better agreement with (i.e., closer than) the
 584 regional reconstructions than this reference. Results are presented for dates when no data are
 585 missing in four periods: 850 to 1350, 1350 to 1850, 850 to 1850, and the full period 850 to
 586 2000. As in section 3.3, we compute the skill in Fig. 4 using the uncertainties provided with
 587 the original reconstruction, as well as a case that assumes the uncertainties are negligible (i.e.,
 588 assuming $\sum e_i^2 = 0$ in Eq. 1 of section 3.4). Additionally, the skill is computed assuming a
 589 reduction by a factor $\sqrt{15}$ in the supplementary figures.

590 The most notable result is that the skill measure is generally undefined when using the
 591 uncertainties provided with the original reconstruction: either the reference or the simulated
 592 data are closer to the reconstruction than uncertainty allows, leading to the square root of a
 593 negative number in Eq. 1. This confirms that uncertainties in the reconstructions are
 594 potentially an overestimation for smoothed time series. When ignoring uncertainties, the 15
 595 year non-overlapping means of the simulations rarely display skill. Simulation skill appears to
 596 be most likely for the European and Arctic regions, while positive skill is nearly absent for the
 597 Southern Hemisphere regions and North America in all the models.

598

599 **4. Links between the different regions**

600 The structure of the spatial variability, i.e. the spatial covariance of temperature
 601 changes, contains contributions from forced signals and from teleconnections in the internal
 602 climate variability. The PAGES 2k temperature reconstructions help to investigate the
 603 consistency between simulations and reconstructions with respect to this covariance structure.
 604 In the following sections, this is evaluated using spatial correlations, Principal Components
 605 (PCs) and Empirical Orthogonal Functions (EOFs), and correlations over sliding temporal
 606 windows.

607 **4.1 Spatial correlation**

608 The spatial correlation matrix of simulated temperature for the PAGES 2k regions is
609 compared to the correlation matrix of the PAGES 2k reconstructions (Fig. 5 and Fig. S6).
610 Correlations are calculated for detrended continental mean time series filtered with a 23 year
611 Hamming window and based on the continents for which these are available, which excludes
612 North America. We use the longest common period for forced simulations and
613 reconstructions, which for the filtered data is 1012 CE – 1978 CE (1000 CE – 1990 CE for
614 annual data). To disentangle the contributions from forcings and from internal variability we
615 analysed forced simulations for the entire analysis period, forced simulations for the pre-
616 industrial period (before 1850 CE), and unforced control simulations.

617 MPI-ESM-P is used to illustrate our main findings in Figure 5 (see Fig. S6 for the other
618 models). Correlations in the forced MPI-ESM-P simulation for the whole period are higher
619 than 0.6 between nearly all regions. In contrast, the correlations for the PAGES 2k
620 temperature reconstructions are rather low, which indicates a substantial inconsistency
621 between the correlation structure in the models and in the PAGES 2k temperature
622 reconstructions. The potential causes of this discrepancy will be discussed in section 6 but we
623 must recall here that, in contrast to other analyses presented above, the evaluation of the
624 spatial correlation does not take into account any uncertainty in the reconstruction. Any non-
625 climatic noise related to the characteristics of the proxy records selected or differences in the
626 reconstruction method between regions would decrease the correlation, contributing to have
627 lower values than for the model results.

628 The correlations in the simulation are lower if only the pre-industrial period is
629 considered, and close to zero in the control simulations. The simulated high correlations for
630 the last century are likely to be a consequence of the rather homogeneous and strong
631 anthropogenic warming in the simulations. The high correlations for the pre-industrial forced
632 runs show that the response to volcanic forcing, solar forcing land use and/or orbital forcing
633 also substantially contributes to the correlations at the time scales considered. Low values
634 obtained for the control simulations indicate that teleconnections between continents are weak
635 for simulated internal variability.

636 Although these general characteristics are present in many of the models evaluated here,
637 there are some differences among them. In particular, some of the models that show higher
638 correlations during pre-industrial times (e.g., CESM) also display a large response to volcanic
639 forcing compared to the other members of the ensemble (Lehner et al., 2015). Additionally,
640 the specific characteristics of some regions may differ substantially. For instance, the
641 correlation between Antarctic temperatures and other regions is very low in MPI-ESM-P or
642 IPSL-CM5A-LR for pre-industrial conditions while it is much larger in CCSM4 and CESM.
643 This can be attributed to a different ratio of forced versus unforced variability, and in
644 particular to discrepancies in the magnitude of the response to external forcing in the selected
645 models.

646 **4.2 Principal Component Analysis**

647 Figure 6a shows the loadings of the first EOF on each region for the PMIP3 forced
648 simulations and the PAGES 2k reconstructions (with corresponding results for the GISS and
649 COSMOS ensembles presented in the supplementary text S4 and Fig. S7). Most models show
650 similarities in the loadings, which indicates that the different regions co-vary similarly in the
651 different models. All loadings are positive and thus the first principal component (PC) is a just
652 weighted mean of all continental temperature series.

653 Consequently, the time series of the first PC of the PMIP3 simulations and PAGES 2k
654 temperature reconstructions (Fig. 6b) reflect the main features of the individual original series
655 (particularly for Northern Hemisphere regions); namely a temperature decline after around
656 1200 CE, which lasts until the early 1800s, followed by the sustained warming within the 19th
657 and 20th century. Additionally, the influence of volcanic eruptions on reconstructed
658 temperatures is visible during some periods, especially during the mid-13th century (although
659 not in the reconstructions), the mid-15th century and the beginning of the 19th century.

660 In most models, the first EOF explains about 80-90% of the total variance, whereas the
661 leading EOF in the PAGES 2k temperature reconstructions accounts for only 55% of the total
662 variance. This shows that the covariance structure is less complex in the simulations. This is
663 consistent with the larger correlations between regions found in section 4.1, which means that
664 the leading mode of homogeneous warming or cooling dominates the covariance structure in
665 model results. In a few simulations (HadCM3, COSMOS), however, the variance explained
666 by the first EOF is about the same as in the reconstructions.

667 The largest values for the loadings are found for the Arctic region, due to the high
668 temperature variability in the last 1,200 years in this region. This expression of the classical
669 Arctic amplification is reflected in most models and in the reconstructions. The ocean-
670 dominated regions of the Southern Hemispheric show less pronounced variability relative to
671 the Northern Hemisphere, consistent with the results of Neukom et al. (2014).

672 If the analysis is performed over the pre-industrial period only (Fig. S8), similar
673 conclusions are reached but the loadings are smaller, especially over the Arctic, and the
674 amount of variance represented in the leading EOF generally decreases, indicating a larger
675 heterogeneity in the pre-industrial period.

676 **4.3 Inter-regional and -hemispheric coherence of past temperature variability**

677 Next, the stationarity of the correlation structure between the different regions, in the
678 models and the reconstructions, is assessed using a running correlation analysis, (Fig. 7, Fig.
679 S9). For the simulations, the multi-model mean shows generally high inter-regional
680 correlations, as the common contribution of the forcing is enhanced because of the averaging
681 procedure. Periods with small variations in external forcing are, however, characterized by
682 weaker coherence between the regions. This occurs during the 11th and 12th century and in
683 shorter periods around 1500 and 1750. High coherence occurs in periods with strong
684 variations in external forcing, highlighting in particular that volcanic eruptions can cause
685 simultaneous temperature variations in most regions.

686 The inter-regional correlations in the individual model simulations vary considerably.
687 The model range includes the correlations derived from the reconstructions for some regions,
688 as for Europe vs. Arctic (Fig. 7a), but values for models are very often higher than for
689 reconstructions (see also section 4.1). The difference is particularly large for the coherence
690 between Australasia and South America (Fig. 7b), which is substantially larger in model
691 simulations compared to reconstructions and instrumental observations (Morice et al., 2012)
692 (Fig. 7b). This could indicate that some regions are less connected by modes of variability
693 (such as ENSO) in reality than suggested by models, that the models have poor representation
694 of modes of internal variability that influence the ocean-dominated Southern Hemisphere (see
695 Neukom et al., 2014; see also Supplementary text S5 and Fig. S10), or that there is more non-
696 temperature noise in the proxy data from those regions.

697

698 **5. Role of Forcing**

699 Some aspects of the response to external forcing have been briefly discussed in the
700 previous sections. It is now formally addressed here by a superposed epoch analysis, by
701 applying the U_R and U_T (correlation- and distance-based) model evaluation statistics and by
702 detection and attribution techniques.

703 **5.1 Superposed epoch analysis**

704 The response of the PAGES 2k reconstructions and the various model simulations to
705 external forcing from solar and volcanic activity is evaluated here using a superposed epoch
706 analysis approach, following Masson-Delmotte et al. (2013). This analysis was conducted for
707 two different time scales, interannual and multidecadal. For interannual time scales, this is
708 done by generating composites of reconstructed and simulated temperature sequences
709 corresponding to the timing of the 12 strongest volcanic events (see supplement sect S2.2).
710 For the multidecadal composites, the 5 strongest events are selected and the means from 40
711 years before to 40 years after the eruption are calculated from time series smoothed with a 40-
712 year low pass filter using least-squares coefficients (Bloomfield, 1976). We also calculate
713 composites corresponding to the timing of intervals of weaker solar forcing at decadal
714 timescales. The intensity of the average model response to the selected forcing events is then
715 compared to the corresponding response found in the reconstructions.

716 The regional temperature response for six PAGES 2k regions (North America is not
717 analysed here; see Section 2.1) to the major volcanic events in the Crowley and Unterman
718 (2012) reconstruction are shown in Figure 8. The temperature perturbation typically lasts
719 longer than the forcing itself, with a recovery to pre-eruption temperatures after 3 to 10 years
720 in the simulations and in the reconstructions.

721 The responses vary considerably in the simulations and in the reconstructions among
722 regions. Nevertheless, the composite averages are always larger in model results with values
723 of up to -1°C compared to about -0.25°C in reconstructions. The largest responses in
724 simulated and reconstructed temperatures are found in Europe and Asia. The Arctic and South
725 America show smaller simulated temperature changes compared to Europe and Asia (around -
726 0.5°C) and the average responses in the reconstructions are even smaller but stay at levels of -
727 0.1 to -0.2°C during several years. For the Antarctic region, both the simulated and
728 reconstructed temperature response is negligible. This is also the case for the reconstructed
729 response in Australasia. Similar results were obtained using the Gao et al. (2008) forcing (see
730 supplementary text S6 and Fig. S11).

731 At multidecadal timescales, the simulated and reconstructed temperature responses are
732 in better agreement, in particular when using the Crowley and Unterman (2012)
733 reconstruction (Fig. S12) rather than the Gao et al. (2008) estimations (Fig. S13), with
734 temperature decreases on the order of a few tenths of degree in most regions. The one
735 exception is South America where, in contrast to simulations, the reconstructions do not show
736 any multidecadal changes associated with volcanic forcing.

737 The multidecadal impact of solar forcing in the reconstructions is strongest in Europe,
738 the Arctic and Asia (Fig. S14), with mean changes ranging from 0.15 to 0.25°C . Changes in
739 model simulations are smaller, lying between 0.05 and 0.1°C in all regions except for

740 Antarctica where no changes are perceptible. The reconstructed changes thus appear larger
741 than the simulated ones in Europe and the Arctic. This interpretation of the results should be
742 approached cautiously, however, as the solar variability is not independent from the volcanic
743 forcing analysis. Volcanic eruptions tend to occur more often in periods of low solar forcing in
744 the reconstructed forcing records, and solar forcing itself is characterised by significant
745 uncertainties (e.g., Schmidt et al., 2011).

746 **5.2 Framework for evaluation of climate model simulations: U_R and U_T statistics**

747 A statistical framework for evaluating simulated temperature sequences against
748 reconstructed past temperature variations was developed by Sundberg et al. (2012), Hind et al.
749 (2012) and Moberg et al. (2015). It includes two essential metrics, which both serve as
750 statistical tests of a null hypothesis. First, a correlation metric, U_R , is used to test whether a
751 significant positive correlation exists between simulated and observed (or reconstructed)
752 temperature variations, indicating that they share a common response to changes in external
753 forcings.

754 Second, a weighted square-distance metric, U_T , is used to test whether temperature
755 variations in a forced simulation are significantly closer to the observed temperature
756 variations than an unforced control simulation. If this is the case, a negative U_T is obtained,
757 whereas a positive U_T indicates that the simulated response to forcings is larger than those in
758 the observations, provided a significant positive U_R is found. Both metrics are approximately
759 distributed as $N(0,1)$ under the null hypothesis that forced simulations are equivalent to
760 unforced control simulations. Thus, it is easy to see if a U_R or U_T value is significantly
761 different from zero. For example, a one-sided test value numerically larger than 1.65 is
762 significant at the 5% level.

763 Prior to the analysis, all records were recalibrated against their instrumental target
764 temperature time series (see supplement sect S2.3) following the procedure of Sundberg et al.
765 (2012) and Moberg et al. (2015) to ensure that each regional reconstruction, after calibration,
766 approximately satisfies the assumption that the true temperature component, upon which the
767 non-climatic noise component is added, is correctly scaled (see Sundberg et al., 2012 and
768 Moberg et al., 2015).

769 Figure 9 shows the model evaluation statistics U_R and U_T (Sundberg et al., 2012),
770 calculated for the 861-1850 pre-industrial period. In general, all forced simulations and the
771 reconstructions share a common forcing signal and, overall the forced simulations match the
772 reconstructions significantly better than the unforced control simulations. However, these
773 overall positive results are essentially due to a good match between simulations and
774 reconstructions in the Northern Hemisphere while the agreement is poorer in the Southern
775 Hemisphere.

776 Because of the imprint of the forcing response, all forced simulations show significant
777 ($p < 0.01$) positive correlation statistics (U_R) when data from all seven regions are combined,
778 although notable differences are seen between regions. In the Arctic, Europe and Asia, all
779 simulations have significant positive U_R values. Nearly all simulations for Australasia and
780 most for Antarctica also have significant positive U_R .

781 In contrast, simulated and reconstructed pre-industrial temperature histories for South
782 America show little common variation, as revealed by mostly insignificant U_R (some are even

783 negative) in that region. U_R statistics for North America (tree-ring based reconstruction) are
784 only slightly better, but note that this reconstruction only starts in 1201. Moreover, the
785 original temporal resolution of 10 years in the North American reconstruction leads to some
786 loss of information in this analysis, which was performed at a 15-year resolution.

787 Results for the distance statistic (U_T) show that nearly all forced simulations are
788 significantly closer ($p < 0.05$) to the observed temperature variations than their respective
789 control simulations when all regions are combined, i.e. their U_T statistics are negative and
790 statistically significant. The Arctic shows the overall best performance in the sense of having
791 the largest number of negative significant U_T values. Most simulations also show negative U_T
792 for Europe, Asia and Antarctica, but many of them are insignificant. For Australasia and South
793 America, nearly all U_T values are insignificant and many are even positive.

794 Thus, overall, the comparison between simulation results and reconstructions performs
795 notably better for the Northern Hemisphere than for the Southern Hemisphere. In particular,
796 nearly all simulations have significant negative U_T values for the combined Northern
797 Hemisphere data ($p < 0.05$) but no significant negative values are found for the Southern
798 Hemisphere where most of the U_T values are positive. This suggests that the simulated effect
799 of forcings in the Northern Hemisphere agrees well in amplitude with the corresponding
800 effect in the Northern Hemisphere reconstructions, whereas the simulated Southern
801 Hemisphere effect of forcings often appears to be larger than in the Southern Hemisphere
802 reconstructions.

803 Results for both U_R and U_T suggest that the most robust agreements are for the largest
804 spatial scales and for ensemble mean results (Fig 9). The most significant U_R and U_T are for
805 ensemble means and global comparisons, followed by ensemble means and NH comparisons.
806 Splitting the analysis period into two halves (856-1350 and 1356-1850, Figs. S15-S16) shows
807 that the more recent period has better U_R statistics. There are, however, not many significant
808 negative U_T in this period, although North America in particular shows several significant
809 values. Extending the analysis to the full 861-2000 period yields higher U_R values for most
810 regions (Fig. S17).

811 The exception is Antarctica, where lower U_R values indicate a divergence of the
812 simulations and reconstruction for this region during the industrial period. Notably, U_T values
813 for the full analysis period are mostly weaker than for the pre-industrial period. Consequently,
814 the overall performance of the simulation results-reconstruction comparison degrades in terms
815 of the distance measure when recent data are included. This is likely because the simulated
816 signal itself often has a larger amplitude in the industrial period than many of the regional
817 temperature reconstructions (see Fig. 1).

818 Ensemble means for COSMOS and GISS ensembles give more significant U_R and U_T
819 than individual simulations from these ensembles, demonstrating once more the value of
820 averaging for isolating the forced signal. The intra-ensemble spread of test statistics illustrates
821 the degree of randomness in U_R and U_T statistics for individual simulations, highlighting the
822 danger of judging one model as being better than another. In particular, it is difficult to judge
823 whether the high (E2) or low (E1) solar forcing amplitude of the COSMOS simulations
824 provides a better fit to the reconstructions, as their ranges of U_T values for individual
825 simulations mostly overlap. For the early period analysis (856-1350), however, the low solar
826 COSMOS simulations provide a better fit than the high solar simulations, as seen by their

827 respective U_T values being of different sign and having entirely non-overlapping ranges when
828 all seven regions or when the Northern Hemisphere regions are combined (Fig. S15).

829 This result is confirmed by a formal test where U_T is calculated in a different way to
830 directly compare the two COSMOS ensembles, using the method described in Moberg et al.
831 (2015). This test reveals that, despite a significantly better fit of the low solar simulations in
832 the earliest period, neither of the two solar forcing alternatives gives a significantly better fit
833 to the reconstructed temperature history when the more recent data are included (Fig. S18).

834 **5.3 Detection and attribution**

835 Detection and attribution aims to identify the forced response in the regional
836 temperature reconstructions by evaluating if observed changes could be entirely caused by
837 variability created within the climate system (internal variability), if external forcing is
838 necessary to explain them, and what magnitude of external forcing response is consistent with
839 reconstructions (see Bindoff et al., 2013; Hegerl and Zwiers, 2011). Here, we focus on all
840 forcings together and not on the response to each forcing individually, as simulations with
841 individual forcings are needed to analyse the latter. Attribution is achieved by estimating the
842 response to the external forcing in the reconstruction using a total least squares regression
843 techniques (following Schurer et al 2013, see also Allen and Stott, 2003). The outcome are
844 scaling factors that determine the amplitude of the fingerprints of the forcing response in the
845 reconstructions. A forcing is detected if a scaling value of zero can be rejected at the 5%
846 significance level, indicating that it is unlikely that climate variability alone is responsible for
847 the similarity between forced response and reconstruction. If the 5-95% range of scaling
848 factors encompasses one then the magnitude of the response to forcing is found to be
849 consistent in simulations and in the reconstruction (see supplement sect. S2.4).

850 Figure 10 shows the results of the detection and attribution analysis using the multi-
851 model ensemble mean, which is calculated as the mean of all model simulations described in
852 section 2.2, except for the high-solar COSMOS and CESM1 simulations as they include a
853 clearly different forcing. All reconstructions and models used were first filtered using a 23-
854 year Hamming window.

855 The response to external forcing is detectable ($p < 0.05$) in all four reconstructions from
856 the Northern Hemisphere and during all time periods (scaling range always greater than zero;
857 indicating that the level of agreement between the multimodel mean and the reconstructions
858 exceeds that from random control samples significantly). The scaling ranges always
859 encompass the scaling factor 1, which shows that the model results are consistent with the
860 reconstructions because they do not need to be scaled up or down.

861 The only exceptions are the earliest time period (864-1350) for North America (tree-ring
862 based reconstruction) where only 150 years of data were available and the early European
863 period, which fails the residual consistency check, indicating that the residual that is attributed
864 to internal variability is larger than expected from model simulations, possibly due to non-
865 climatic noise in reconstructions. The results for the latter case suggest that the basic
866 hypotheses underlying the methodology are violated because the model-reconstruction
867 discrepancy cannot be explained by internal variability alone. External forcing is also
868 detectable when the model and reconstruction data from all Northern Hemisphere regions are
869 combined.

870 External forcing is not detectable in South America (no scaling ranges significantly
871 larger than 0) and only for certain time periods for Antarctica and Australasia (with fits for
872 Australasia failing the residual consistency check). External forcing is also not detectable in
873 the combined Southern Hemisphere reconstruction. As well as being un-detectable, despite
874 accounting for uncertainty in simulated signals due to variability, the estimated signals are
875 also significantly smaller than simulated. Consequently, the models appear to simulate a
876 magnitude and pattern of external forcing in the Southern Hemisphere significantly different
877 from that derived from the PAGES 2k reconstructions. This could be due to strong noise in
878 reconstructions swamping the forced response, calibration uncertainty in reconstructions
879 misestimating the magnitude of the forced response, or errors in climate models as discussed
880 below

881

882 **6. Discussion**

883 In the light of the results presented in the Sections 3 to 5, we discuss below each of the
884 three questions raised in the introduction.

885 **6.1 Are the statistical properties of surface temperature data for each individual** 886 **continent-scale region consistent between simulations and reconstructions?**

887 The analyses herein show that the answer to this question depends on the specific
888 feature assessed. The simulation results and reconstructions agree at regional scale for some
889 metrics, but disagree in many cases. The consistency between simulations and observations is
890 still generally more robust at hemispheric and global scales, and the fit to reconstructions is
891 improved for ensemble mean of simulations compared to individual members.

892 Overall, smoothed simulated temperature anomalies from the long-term average lie
893 within the range of the originally published uncertainty estimates of the reconstructions.
894 However, these uncertainty ranges are, for all regions except North America, defined for data
895 at annual resolution and therefore are very likely larger than uncertainty ranges adapted for
896 the smoothed versions of the data (see section 3.1). Thus, the published uncertainties for the
897 reconstructions are in most cases too large to provide strong constraints on the ensemble of
898 simulations, as different forcing amplitudes and responses are nevertheless consistent within
899 the range of the reconstructed values. Some common signals between model results and
900 reconstructions can be identified visually as isolated events, such as the cooling during the
901 early 19th century in many regions, but they are relatively rare.

902 The time series for forced simulations are nevertheless significantly correlated with
903 temperature reconstructions, for many regions, when the entire series are considered
904 (Supplementary text S3). Models also have some skill compared to a simple *a priori* estimate
905 assuming no temperature change over the past millennium (section 3.4). Despite using a very
906 simple reference method as a benchmark, however, such skill is achieved nearly exclusively
907 for Northern Hemisphere regions, specifically for the Arctic and in some models for Europe
908 and Asia. This is in agreement with the conclusions derived from the application of the
909 Sundberg et al. (2012) evaluation framework (section 5.2) that forced simulations are
910 significantly closer to the reconstructions than unforced simulations in Northern Hemisphere
911 regions. In particular, the Arctic region shows a robust agreement, as do Europe and Asia to a
912 lesser degree. In contrast, for all the regions of the Southern Hemisphere, the models have

913 nearly no skill compared to a constant climate reference and individual forced simulations are
914 in most cases not significantly closer to reconstructions than an unforced reference.

915 The diagnostics mentioned above addressed whether simulated time series of surface
916 temperature at continental scale have temporal similarities with reconstructed ones. The
917 climatological or probabilistic consistency is complementary as it focuses on the distribution
918 of temperature data, independent of the particular trajectory over time. For most regions, no
919 consistency is found between the distribution of model results and reconstructed temperatures
920 when using the original reconstruction uncertainty estimates (section 3.3), which are annually
921 resolved in all cases except North America.

922 It should be noted, however, these results depend strongly on the uncertainty estimates
923 considered (Bothe et al., 2013a,b): the greater the assumed reconstruction uncertainty, the
924 weaker the consistency with model simulations as the models tend to appear over-dispersive.
925 When reducing the uncertainties, to adapt them to the smoothing or temporal averaging
926 applied here, the consistency improves in many regions. Such reduction of the uncertainties
927 may, however, lead to overconfident conclusions if the original uncertainty estimates at the
928 annual resolution did not account for all existing sources of uncertainty.

929 A visual comparison suggests that the temperature reconstructions show slightly more
930 centennial to multi-centennial variability over the full period with stronger long-term trends,
931 while model results indicate a stronger recent warming compared to some of the
932 reconstructions (section 3.1). Comparison of the series spectra (section 3.2) reveals marked
933 differences between regions in how well the simulations agree with the reconstructions. The
934 best overall agreement is seen for the Arctic, where the model spectra mostly lie within a 95%
935 confidence interval for the reconstruction spectrum. For all other regions, the model spectra
936 often lie outside the confidence interval for some frequency ranges. The mismatch is most
937 pronounced for South America, but there are other examples with both lower and higher
938 variance at different frequencies in model results compared to reconstructions.

939 The disagreements can have various origins, in either reconstructions or simulations or
940 both. For instance, the total variance of reconstructions is dependent on how they were
941 calibrated to instrumental observations (e.g., Kutzbach et al., 2011) but the shape and slopes
942 of their spectra are determined by spectra of both the true climate and the non-climatic proxy-
943 data noise and by the signal-to-noise ratio (Moberg et al., 2008). Some studies have suggested
944 that reconstruction methodologies may alone underestimate low-frequency variability, in
945 addition to any frequency biases inherent to the proxy data (e.g., Smerdon et al., 2010; Esper
946 et al., 2012; Smerdon et al., 2015). The amplitude of the reconstructed past forcing changes,
947 which affect the model spectra, is still uncertain (Schmidt et al., 2011, 2012). The modelled
948 transient climate response and the amplitude of internal variability at the regional scale vary
949 considerably and thus deficiencies in applied forcings or internal model physics can lead to
950 errors in the modelled spectra. Nevertheless, no major, systematic model underestimation of
951 low frequency variability can be deduced at the continental scale from the analyses performed
952 herein, in contrast to some recent studies devoted to the ocean surface temperature (Laepple
953 and Huybers, 2014ab).

954 **6.2 Are the cross regional relations of temperature variations similar in**
955 **reconstructions and models?**

956 Discrepancies in the interregional relations between reconstructions and model results
957 are clearer than for each individual region. While the strong correlations between the
958 temperature variations in regions from the Northern Hemisphere in model simulations have
959 some similarities to the ones in the reconstructions (section 4.1), the correlation between the
960 hemispheres and between the Southern Hemisphere regions are much stronger in models than
961 in reconstructions, as previously reported by Neukom et al. (2014) at hemispheric scale.

962 This result is robust as it is also reflected in the larger variance explained by the first
963 EOF mode in models than in the temperature reconstructions (section 4.2) and this is valid for
964 most of the past millennium (section 4.3). These differences may be due to a stronger
965 response to forcings in the models, to unrealistic internal variability in the models, or to non-
966 climatic noise in the proxies or due to a combination of these factors, as discussed in more
967 detail below. Additionally, there are large differences between the various models in the
968 Southern Hemisphere. For instance Antarctic temperature is strongly related to other regional
969 temperatures in some simulations and not in others, suggesting that specific model dynamics
970 may account for some of the discrepancies with the reconstructions.

971 **6.3 Can the signal of the response to external forcing be detected on continental** 972 **scale and, if so, how large are these signals?**

973 The agreements or disagreements between model results and reconstructions can be
974 partly explained by the model response to forcing. The contribution of the forcing derived
975 from simulated results can be detected in the reconstructions for all regions of the Northern
976 Hemisphere (section 5.3). The forcings used in the PMIP3 model experiment result in
977 simulated temperature histories that, on the whole, explain a significant fraction of the past
978 regional temperatures in the pre-industrial climate.

979 This strongly contributes to the model skill for the Northern Hemisphere, as unforced
980 internal stochastic variability is unlikely to agree between model results and observations.
981 This is confirmed by the significant correlation coefficients (Fig. S1) and correlation statistics
982 (U_R) (section 5.2) that indicate common external forcing variations. Furthermore, the
983 correlations increase for the ensemble average relative to the available single-model
984 simulations due to the fact that the contributions from internal variability are reduced by
985 averaging.

986 On interannual time scales, the model response to volcanic forcing appears larger than
987 represented in the reconstructions (section 5.1). There is some debate on the potential
988 underestimation or overestimation of the cooling due to volcanic eruptions in reconstructions
989 (e.g., Mann et al., 2012; Anchukaitis et al., 2012; Tingley et al., 2014; Büntgen et al., 2015).
990 Nevertheless, this model overestimation was also found when compared to instrumental data
991 and at hemispheric scale, suggesting a robust phenomenon (Brohan et al., 2012; Fernandez-
992 Donado et al., 2013; Masson-Delmotte et al., 2013; Schurer et al., 2013). Both model results
993 and reconstructions also show that volcanic activity impacts temperature at multidecadal
994 timescales, with a similar magnitude of the temperature response in models and
995 reconstructions over most of the regions in the Northern Hemisphere. This is consistent with
996 the detection and attribution analysis (section 5.3), which indicates that the magnitude of the
997 simulated response to forcing in the Northern Hemisphere has the correct amplitude for
998 smoothed time series.

999 The role of solar forcing is less clear and none of the pre-PMIP3 COSMOS simulations
1000 with either a moderate or a weak solar forcing gives a systematically better agreement with
1001 the reconstructions than the other, although the ensemble with low solar forcing yield a better
1002 fit during the first 500 years (Fig. S18). This confirms earlier results obtained at the
1003 hemispheric scale (Masson-Delmotte et al., 2013; Schurer et al., 2014).

1004 In the Southern Hemisphere, the influence of external forcing is often not detected
1005 (section 5.3). This is consistent with the lower correlation coefficients (Fig. S1) and weaker
1006 correlation statistics (U_R) there (section 5.2). The models also seem to overestimate the
1007 response compared to the signal recorded in the Southern Hemisphere reconstructions (section
1008 5.2-5.3). This finding is likely related to the larger covariability seen within Southern
1009 Hemisphere regions in models compared to reconstructions. Moreover, control simulations
1010 display low correlations between the Northern and Southern Hemisphere regions.

1011 The analysis performed herein, however, cannot reveal the origin of the mismatch
1012 between simulation results and reconstructions. These differences may be due to biases in the
1013 dynamics of the climate models or to errors in the implemented forcing, in particular in their
1014 spatial distribution. Land-use changes, which are not included in some models (Table 1), tend
1015 to reduce the spatial correlation between regions as deforestation did not occur at the same
1016 time over all continents (Pongratz et al., 2008; Kaplan et al., 2011).

1017 The spatial distribution of volcanic aerosols may also contribute to pronounced regional
1018 differences. Volcanic forcing is usually not implemented as a direct simulation of changes in
1019 stratospheric sulphate concentrations due to individual eruptions, but as a mean change in the
1020 optical depth for different latitudinal bands. This can have an impact on the overestimation of
1021 the response in individual simulations or to individual eruptions. Additionally, if the
1022 latitudinal distribution of volcanic aerosols is too homogeneous, thereby inducing
1023 unrealistically symmetric forcing between hemispheres, it would also overestimate the global
1024 signature of the induced cooling.

1025 Any non-climatic noise in the reconstruction would tend to reduce the covariance in
1026 reconstructions compared to model results, which would lead to an underestimation of the
1027 relative contribution of the forced signal. Despite the large progress made over the last few
1028 years, this may still be a critical problem in the Southern Hemisphere, where fewer long
1029 paleoclimate records are available compared to the Northern Hemisphere, explaining some of
1030 the model-data mismatch there.

1031 The role of internal variability in driving temperature variations may also be
1032 underestimated in model simulations, particularly in the ocean dominated Southern
1033 Hemisphere, as suggested by Neukom et al. (2014). Simulated internal variability may,
1034 however, be overestimated, as reported here in at least one model and elsewhere for ENSO-
1035 type variability (Jungclaus et al., 2006) or for the Southern Ocean ice extent (Zunz et al.,
1036 2013). This would imply a ratio between internal and forced variability that is incorrect,
1037 which would lead to biased correlations between the different regions.

1038 Another potential explanation for the differences in the spatial covariance structure
1039 between models and observations relates to the relatively coarse resolution of the climate
1040 models. Using models with higher spatial resolution will increase the number of spatial
1041 degrees of freedom and potentially improve the co-variance structure of the climate models

1042 compared to reconstructions. Nevertheless, the expense required for both high spatial and
1043 temporal resolution, as well as the necessary ensemble approach could be prohibitive.

1044

1045 **7. Conclusions and perspectives**

1046 The analysis of model simulations and PAGES 2k temperature reconstructions has
1047 allowed us to extend some of the some conclusions previously articulated at only hemispheric
1048 scale. For all the continental-scale regions in the Northern Hemisphere, the models are able to
1049 simulate a forced response with a magnitude similar to the one derived from reconstructions.
1050 Despite higher levels of variability on continental scales (relative to full hemispheres), the
1051 role of forcing is found to be important. This leads to reasonable agreement between models
1052 and temperature reconstructions.

1053 Nevertheless, a deeper assessment of the consistency between simulated results and
1054 reconstruction is limited because of the large uncertainties in the reconstructions and the weak
1055 constraints on the estimates of this uncertainty. Notably, the agreement between simulation
1056 results and reconstructions is poor for the Southern Hemisphere regions. Our results indicate
1057 that models have a much clearer response to forcing than deduced from the reconstructions,
1058 leading to a greater consistency across the Southern Hemisphere regions and between
1059 hemispheres in model results than in the reconstructions.

1060 It is not possible to precisely assess which part of those disagreements comes from the
1061 biases in model dynamics, the forcing or in the reconstructions. As suggested in many
1062 previous studies, substantial progress will only be possible with better uncertainty
1063 quantification and reduction (spatially and temporally) in the reconstructions and the forcing,
1064 and through model improvements.

1065 Nevertheless, on the basis of our results we highlight four specific points that may lead
1066 to significant advances in the coming years.

1067 The first is the insights that can be gained through studying the discrepancies between
1068 reconstructions and simulations relative to direct observations over the most recent decades. A
1069 quantitative comparison between simulations, reconstructions, and instrumental data would
1070 provide useful information on the origin of those disagreements, allow an estimate of the non-
1071 climatic noise in reconstructions, and would elucidate how mismatches over the last 150 years
1072 are related to disagreements over the last several millennia (e.g., Ding et al 2014).

1073 Secondly, large uncertainties are associated with the behaviour of the ocean over the
1074 past millennium. The discrepancies in the low frequency variability between model results
1075 and reconstructions at the continental scale seem less systematic than for some oceanic data
1076 (Laepplé and Huybers, 2014ab), but clearly assessing this would require additional analyses.
1077 As new paleoclimate data compilations are now available for the global ocean (Tierney et al.,
1078 2015; McGregor et al., 2015), model-data comparison for oceanic regions should be
1079 encouraged, and the compatibility between ocean and land temperature reconstructions tested.
1080 This would allow us to assess the multidecadal internal and forced variability of the ocean and
1081 to determine if it is the origin of the disagreement between model simulations and Southern
1082 Hemisphere reconstructions (e.g., Neukom et al., 2014). Internal ocean variability can also
1083 have a significant influence on Northern Hemisphere climate as seen in several studies
1084 investigating the circulation in the Atlantic at multi-decadal time scales (e.g., Delworth and

1085 Mann, 2000; Knight et al., 2005; Lohmann et al., 2014). These are the timescales for which
1086 most models tend to display less variability than reconstructions.

1087 Third, our comparison of continental-scale temperature reconstructions with simulated
1088 temperatures only uses a small fraction of the information provided by models and
1089 paleoclimate records. As discussed in Phipps et al. (2013), other approaches can provide
1090 analyses complementary to classical model-data comparison, through a better handling of the
1091 various sources of uncertainty. Promising examples are proxy forward models, which simulate
1092 directly the proxy records from climate model outputs (e.g., Evans et al., 2013) and data
1093 assimilation methods (e.g., Widmann et al., 2010; Goosse et al., 2012b; Steiger et al., 2014).
1094 These approaches combine model results and observations to obtain the best estimates of past
1095 change and may be most effective at to detecting inconsistencies between model and
1096 palaeoclimate estimates.

1097 Finally, one could also question the selection of the continental scale as the basis of a
1098 comparison, as regional changes are strongly affected by modes of variability such as ENSO,
1099 the Southern Annular Mode, the North Atlantic Oscillation or the Pacific North America
1100 pattern. These modes could imprint temperature patterns that are masked by averaging over
1101 the continents. On the other hand, model-data comparison made at smaller spatial scales has
1102 revealed highly variable and even contradictory results at nearby regions (Moberg et al.,
1103 2015), suggesting that a large number of local proxy data sites are needed for obtaining robust
1104 results. Ideally, a sub-regional selection from key teleconnection regions should be used to
1105 assess the stability of climate modes (Raible et al., 2014) or enable reliable reconstruction of
1106 modes of variability (Lehner et al., 2012; Zanchettin et al., 2015; Ortega et al., 2015),
1107 although this requires strong reconstruction skill to be successful (e.g., Russon et al., 2015).
1108 Additionally, spatially resolved reconstructions should be targeted because they offer useful
1109 potential for dynamic interpretation (e.g., Luterbacher et al., 2004; Mann et al., 2009; Steiger
1110 et al., 2014, PAGES 2k Consortium, 2014; Shi et al., 2015).

1111 In summary, our results for the Northern Hemisphere suggest a convergence of our
1112 understanding of climate variability over the past 1000 years, but there remain many open
1113 questions for the Southern Hemisphere. Progress may be expected from comparing
1114 simulations, reconstructions and observations in the instrumental period, from a better
1115 knowledge of internal and forced variability of the ocean, from efforts to understand climate
1116 variability via proxy forward modelling and data assimilation, and from a clearer view of the
1117 influence of climate modes on temperature variability.

1118

1119 **Acknowledgements**

1120 This study is based on discussions held during the joint workshop of the PAGES 2k network
1121 and PAST2k-PMIP *Integrated analyses of reconstructions and multi-model simulations for*
1122 *the past two millennia*, Madrid, Spain, 4-6 November 2013. PAGES and FECYT (FCT-13-
1123 6276) are greatly thanked for supporting this workshop. We acknowledge the World Climate
1124 Research Programme's Working Group on Coupled Modelling, which is responsible for
1125 CMIP. The U.S. Department of Energy's Program for Climate Model Diagnosis and
1126 Intercomparison provides coordinating support for CMIP and led the development of software
1127 infrastructure in partnership with the Global Organization for Earth System Science Portals.
1128 H.G. is Research Director with the Fonds National de la Recherche Scientifique (F.R.S.-
1129 FNRS-Belgium). This work is supported by the F.R.S.- FNRS and by the Belgian Federal

1130 Science Policy Office (Research Program on Science for a Sustainable Development). CCR
1131 and FL are supported by the Swiss National Science foundation. PY is supported by the
1132 MILEX project of the Swedish Research Council. JG is funded by Australian Research
1133 Council Project DE130100668. OB was supported by LOCHMES (Leibniz-Society), PRIME-
1134 II (within DFG INTERDYNAMIK) and CliSAP. LFD was funded by a FPU grant: AP2009-
1135 4061. AM and AH are supported by the Swedish Research Council grants B0334901 and
1136 C0592401. GH and AS are supported by the ERC advanced grant TITAN (320691).

1137 **Author contributions**

1138 LFD, FGR, EGB, HG, JJ organized the workshop at the origin of this paper. HG led the
1139 synthesis. OB, HG, GH, AM, CR, AS, SW, EZ coordinated the writing. LFD, EGB AH, FL,
1140 NM prepare the data sets and made them available to the whole group. OB, LFD, EGB, FGR,
1141 AH, FL, NM, AM, AS, SW, EW, MW, EZ performed the figures and their initial analysis. All
1142 authors contribute to the writing of the various sections and reviewed the manuscript.

1143

1144 References

- 1145 Abram, N. J., Mulvaney, R., Vimeux, F., Phipps, S. J., Turner, J. and England, M. H.:
1146 Evolution of the Southern Annular Mode during the past millennium, *Nat. Climate Change*, 4,
1147 564–569, 2014.
- 1148 Allen, M. R., and P. A. Stott: Estimating signal amplitudes in optimal fingerprinting, Part I:
1149 Theory, *Clim. Dyn.* 21, 477-491, 2003.
- 1150 Ammann, C. M., Joos, F., Schimel, D. S., Otto-Bliesner, B. L., and Tomas, R. A.: Solar
1151 influence on climate during the past millennium: Results from transient simulations with the
1152 NCAR Climate System Model. *Proc. Nat. Acad. Sciences*, 104(10), 3713-3718. Doi:
1153 10.1073/pnas.0605064103, 2007.
- 1154 Anchukaitis, K.J., Breitenmoser, P., Briffa, K. R., Buchwal, A., Büntgen, U., Cook, E. R.,
1155 D'Arrigo, R. D., Esper, J., Evans, M. N., Frank, D., Grudd, H., Gunnarson, B. E., Hughes, M.
1156 K., Kirilyanov, A. V., Körner, C., Krusic, P. J., Luckman, B., Melvin, T. M., Salzer, M. W.,
1157 Shashkin, A.V., Timmreck, C., Vaganov, E. A., and Wilson, R. J. S.: Tree rings and volcanic
1158 cooling, *Nat. Geosci.*, 5(12), 836–837, doi:10.1038/ngeo1645, 2012.
- 1159 Anderson, J. L.: A method for producing and evaluating probabilistic forecasts from ensemble
1160 model integrations, *J. Climate*, 9, 1518–1530, 1996
- 1161 Annan, J. D., and Hargreaves, J. C.: Reliability of the CMIP3 ensemble, *Geophys. Res. Lett.*,
1162 37, L02703, doi:10.1029/2009GL041994, 2010.
- 1163 Ault, T. R., Deser, C., Newman, M. and Emile-Geay, J.: Characterizing decadal to centennial
1164 variability in the equatorial Pacific during the last millennium. *Geophys. Res. Lett.* 40 (13),
1165 3450-3456, doi: 10.1002/grl.50647, 2013.
- 1166 Bard, E., Raisbeck, G., Yiou, F., and Jouzel, J.: Solar irradiance during the last 1200 years
1167 based on cosmogenic nuclides. *Tellus B*, 52(3), 985-992. Doi: 10.1034/j.1600-0889.2000.d01-
1168 7.x, 2000.
- 1169 Berger, A.: Long-term variations of daily insolation and Quaternary climatic changes, *J.*
1170 *Atmos. Sciences*, 35(12), 2362-2367. Doi: 10.1175/1520-
1171 0469(1978)035<2362:LTVODI>2.0.CO;2, 1978.
- 1172 Bertrand, C., Loutre, M.-F., Crucifix, M., and Berger, A.: Climate of the last millennium: a
1173 sensitivity study. *Tellus* 54A, 221-244, 2002.
- 1174 Bindoff, N.L., Stott, P.A., AchutaRao, K.M., Allen, M.R., Gillett, N., Gutzler, D., Hansingo,
1175 K., Hegerl, G., Hu, Y., Jain, S., Mokhov, I.I., Overland, J., Perlwitz, J., Sebbari, R., and Zhang,
1176 X. : Detection and Attribution of Climate Change: from Global to Regional. In: *Climate*
1177 *Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth*
1178 *Assessment Report of the Intergovernmental Panel on Climate Change* [Stocker, T.F., D. Qin,
1179 G.-K. Plattner, M. Tignor, S.K. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex and P.M.
1180 Midgley (eds.)]. Cambridge University Press, Cambridge, United Kingdom and New York,
1181 NY, USA., 2013.
- 1182 Bloomfield, P.: *Fourier Analysis of Time Series: An Introduction*. New York: John Wiley and
1183 Sons, 1976.

- 1184 Borlace, S., Cai, W. and Santoso, A.: Multidecadal ENSO Amplitude Variability in a 1000-yr
 1185 Simulation of a Coupled Global Climate Model: Implications for Observed ENSO Variability,
 1186 *J. Clim.* 26 (23), 9399-9407, 2013.
- 1187 Bothe, O., Jungclaus, J. H., Zanchettin, D., and Zorita, E.: Climate of the last millennium:
 1188 ensemble consistency of simulations and reconstructions, *Clim. Past*, 9, 1089-1110,
 1189 doi:10.5194/cp-9-1089-2013, 2013a.
- 1190 Bothe, O., Jungclaus, J. H., and Zanchettin, D.: Consistency of the multi-model
 1191 CMIP5/PMIP3-past1000 ensemble, *Clim. Past*, 9, 2471-2487, doi:10.5194/cp-9-2471-2013,
 1192 2013b.
- 1193 Braconnot, P., Harrison, S. P., Kageyama, M., Bartlein, P. J., Masson-Delmotte, V., Abe
 1194 Ouchi, A., Otto-Bliesner, B., and Zhao, Y.: Evaluation of climate models using paleoclimate
 1195 data, *Nature Climate Change*, 2, 417–424 doi:10.1038/NCLIMATE1456, 2012.
- 1196 Bretagnon, P., and Francou, G.: Planetary theories in rectangular and spherical variables-
 1197 VSOP 87 solutions. *Astronomy and Astrophysics*, 202, 309-315, 1988.
- 1198 Brohan, P., Allan, R., Freeman, E., Wheeler, D., Wilkinson, C., Williamson, F.: Constraining
 1199 the temperature history of the past millennium using early instrumental observations, *Clim.*
 1200 *Past*, 8, 1551-1563, doi:10.5194/cp-8-1551-2012, 2012
- 1201 Büntgen, U., Trnka, M., Krusic, P. J., Kyncl, T., Kyncl, J., Luterbacher, J., Zorita, E.,
 1202 Charpentier Ljungqvist, F., Auer, I., Konter, O., Schneider, L., Tegel, W., Štěpánek, P.,
 1203 Brönnimann, S., Hellmann, L., Nievergelt, D., and Esper, J.: Tree-Ring Amplification of the
 1204 Early-19th Century Summer Cooling in Central Europe, *J. Clim.*, in press, doi:
 1205 <http://dx.doi.org/10.1175/JCLI-D-14-00673.1>, 2015
- 1206 Christiansen, B., Schmith, T., and Thejll, P.: A surrogate ensemble study of climate
 1207 reconstruction methods: Stochasticity and robustness, *J. Clim.*, 22, 951–976,
 1208 doi:10.1175/2008JCLI2301.1, 2009.
- 1209 Coats, S., Cook, B.I., Smerdon, J.E., and Seager, R.: North American Pan-Continental
 1210 droughts in model Simulations of the last millennium, *J. Clim* 28, 2025-2043, 2015a.
- 1211 Coats, S., Smerdon, J.E., Cook, B.I., and R. Seager, R.: Are simulated megadroughts in the
 1212 North American Southwest forced?, *J. Clim.*, 28, 124-142. doi:[http://dx.doi.org/10.1175/JCLI-](http://dx.doi.org/10.1175/JCLI-D-14-00071)
 1213 [D-14-00071](http://dx.doi.org/10.1175/JCLI-D-14-00071), 2015b.
- 1214 Comboul, M., Emile-Geay, J., Hakim, G. J. and Evans, M. N.: Paleoclimatic sampling as a
 1215 sensor placement problem. *J. Clim.*. 28, 7717-7740, 2015
- 1216 Cook, B.I., Ault, T.R., and Smerdon, J.E.: Unprecedented 21st-century drought risk in the
 1217 American Southwest and Central Plains, *Science Advances*, 1, e1400082, 2015.
- 1218 Cook, E.R., Meko, D. M., Stahle, D. W., and Cleaveland, M. K: Drought reconstructions for
 1219 the continental United States. *J. Climate*, 12, 1145–1162. doi: [http://dx.doi.org/10.1175/1520-](http://dx.doi.org/10.1175/1520-0442(1999)012<1145:DRFTCU>2.0.CO;2)
 1220 [0442\(1999\)012<1145:DRFTCU>2.0.CO;2](http://dx.doi.org/10.1175/1520-0442(1999)012<1145:DRFTCU>2.0.CO;2)., 1999.
- 1221 Cook, E.R., Woodhouse, C. A., Eakin, C. M., Meko, D. M., Stahle, D. W.: Long-term aridity
 1222 changes in the Western United States, *Science* 306, 1015, doi: 10.1126/science.1102586,
 1223 2004.

- 1224 Cook, E.R., Anchukaitis, K. J., Buckley, B. M., D'Arrigo, R. D., Jacoby, G. C., Wright, W. E.:
 1225 Asian monsoon failure and megadrought during the Last Millennium, *Science*, 328, 486,
 1226 2010.
- 1227 Crowley, T.J.: Causes of Climate Change Over the Past 1000 Years, *Science*, 289, 270-277,
 1228 2000.
- 1229 Crowley, T.J. and Unterman, M. B.: Technical details concerning development of a 1200-yr
 1230 proxy index for global volcanism, *Earth System Science Data*, Vol. 5, pp. 187-197, DOI:
 1231 10.5194/essd-5-187-2013, 2013.
- 1232 D'Arrigo, R., Wilson, R., Jacoby, G., On the long-term context for late twentieth century
 1233 warming, *J. Geophys. Res.*, 111, D03103, doi:10.1029/2005JD006352, 2006.
- 1234 Delworth, T.L., and Mann, M.E.: Observed and simulated multidecadal variability in the
 1235 Northern Hemisphere, *Clim. Dyn* 16:661-676, 2000.
- 1236 Ding, Y., Carton, J. A., Chepurin, G. A., Stenchikov, G., Robock, A., Sentman, L. T., and
 1237 Krasting, J. P.: Ocean response to volcanic eruptions in Coupled Model Intercomparison
 1238 Project 5 simulations, *J. Geophys. Res. Oceans*, 119, 5622–5637, doi:10.1002/2013JC009780,
 1239 2014.
- 1240 Dufresne, J.-L., Foujols, M.-A., Denvil, S., Caubel, A., Marti, O., Aumont, O., Balkanski, Y.,
 1241 Bekki, S., Bellenger, H., Benshila, R., Bony, S., Bopp, L., Braconnot, P., Brockmann, P.,
 1242 Cadule, P., Cheruy, F., Codron, F., Cozic, A., Cugnet, D., de Noblet, N., Duvel, J.-P., Ethé, C.,
 1243 Fairhead, L., Fichefet, T., Flavoni, S., Friedlingstein, P., Grandpeix, J.-Y., Guez, L., Guilyardi,
 1244 E., Hauglustaine, D., Hourdin, F., Idelkadi, A., Ghattas, J., Joussaume, S., Kageyama, M.,
 1245 Krinner, G., Labetoulle, S., Lahellec, A., Lefebvre, M.-P., Lefevre, F., Levy, C., Li, Z. X.,
 1246 Lloyd, J., Lott, F., Madec, G., Mancip, M., Marchand, M., Masson, S., Meurdesoif, Y.,
 1247 Mignot, J., Musat, I., Parouty, S., Polcher, J., Rio, C., Schulz, M., Swingedouw, D., Szopa, S.,
 1248 Talandier, C., Terray, P., Viovy, and N., Vuichard, N.: Climate change projections using the
 1249 IPSL-CM5 Earth System Model: from CMIP3 to CMIP5, *Clim. Dyn.* 40(9-10), 2123-2165,
 1250 2013.
- 1251 Emile-Geay, J., Cobb, K. M., Mann, M. E., and Wittenberg, A. T.: Estimating Central
 1252 Equatorial Pacific SST Variability over the Past Millennium. Part I: Methodology and
 1253 Validation, *J. Climate*, 26, 2302–2328, doi: <http://dx.doi.org/10.1175/JCLI-D-11-00510.1>,
 1254 2013
- 1255 Esper, J., Frank, D., Wilson, R., and Briffa, K.: Effect of scaling and regression on
 1256 reconstructed temperature amplitude for the past millennium, *Geophys. Res. Lett.*, 32,
 1257 L07711, doi:10.1029/2004GL021236, 2005.
- 1258 Esper, J., Frank, D.C., Timonen, M., Zorita, E., Wilson, R. J. S., Luterbacher, J., Holzkämper,
 1259 S., Fischer, N., Wagner, S., Nievergelt, D., Verstege, A, and Büntgen, U.: Orbital forcing of
 1260 tree-ring data, *Nature Clim. Change*, 2, 862–866, 2012.
- 1261 Evans, M.N., Tolwinski-Ward, S.E., Thompson, D.M., and Anchukaitis, K.J.: Applications of
 1262 proxy system modeling in high resolution paleoclimatology, *Quate. Science Rev.* 76, 16-28,
 1263 2013.

- 1264 Fernández-Donado, L., González-Rouco, J. F., Raible, C. C., Ammann, C. M., Barriopedro,
 1265 D., García-Bustamante, E., Jungclaus, J. H., Lorenz, S. J., Luterbacher, J., Phipps, S. J.,
 1266 Servonnat, J., Swingedouw, D., Tett, S. F. B., Wagner, S., Yiou, P., and Zorita, E.: Large-scale
 1267 t temperature response to external forcing in simulations and reconstructions of the last
 1268 millennium, *Clim. Past*, 9, 393-421, DOI: [doi:10.5194/cp-9-393-2013](https://doi.org/10.5194/cp-9-393-2013), 2013.
- 1269 Feulner, G.: Are the most recent estimates for Maunder Minimum solar irradiance in
 1270 agreement with temperature reconstructions? *Geophysical Research Letters*, 38, L16706,
 1271 doi:10.1029/2011GL048529, 2011.
- 1272 Flato, G., Marotzke, J., Abiodun, B., Braconnot, P., Chou, S.C., Collins, W., Cox, P., Driouech,
 1273 F., Emori, S., Eyring, V., Forest, C., Gleckler, P., Guilyardi, E., Jakob, C., Kattsov, V., Reason,
 1274 C. and Rummukainen, M.: Evaluation of Climate Models. In: *Climate Change 2013: The
 1275 Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of
 1276 the Intergovernmental Panel on Climate Change* [Stocker, T.F., D. Qin, G.-K. Plattner, M.
 1277 Tignor, S.K. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex and P.M. Midgley (eds.)].
 1278 Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA, 2013.
- 1279 Flückiger, J., Dällenbach, A., Blunier, T., Stauffer, B., Stocker, T. F., Raynaud, D., and
 1280 Barnola, J. M.: Variations in atmospheric N₂O concentration during abrupt climatic changes.
 1281 *Science*, 285(5425), 227-230. doi: 10.1126/science.285.5425.227, 1999.
- 1282 Flückiger, J., Monnin, E., Stauffer, B., Schwander, J., Stocker, T. F., Chappellaz, J., and
 1283 Barnola, J. M.: High resolution Holocene N₂O ice core record and its relationship with CH₄
 1284 and CO₂, *Glob. Biogeochem. Cycles*, 16(1), 10-1, doi: 10.1029/2001GB001417, 2002
- 1285 Franke, J., Frank, D., Raible, C. C., Esper, J. and Brönnimann, S.: Spectral biases in tree-ring
 1286 climate proxies, *Nat. Clim. Change*, 3, 360–364, doi:10.1038/nclimate1816, 2013.
- 1287 Gao, C., Robock, A., and Ammann, C.: Volcanic forcing of climate over the last 1500 years:
 1288 An improved ice-core based index for climate models, *J. Geophys. Res.*, 113, D2311,
 1289 doi:10.1029/2008JD010239, 2008.
- 1290 Gergis, J., Neukom, R., Gallant, A. J. E. and Karoly, D. J.: Evidence of rapid late 20th century
 1291 warming from Australasian temperature reconstruction ensembles spanning the last
 1292 millennium, *J. Clim.* (submitted), 2015.
- 1293 Gonzalez-Rouco, J.F., Beltrami, H., Zorita, E., and von Storch, H.: Simulation and inversion
 1294 of borehole temperature profiles in surrogate climates: spatial distribution and surface
 1295 coupling, *Geophys. Res. Lett.* 33, L01703, doi:10.1029/2005GL024693, 2006.
- 1296 Goosse, H., Masson-Delmotte, V., Renssen, H., Delmotte, M., Fichefet, T., Morgan, V., van
 1297 Ommen, T., Khim, B. K., and Stenni, B.: A late medieval warm period in the Southern Ocean
 1298 as a delayed response to external forcing? *Geophysical Research Letters* 31 (6): L06203,
 1299 2004.
- 1300 Goosse, H., Renssen, H., Timmermann, A., and Bradley, R.S.: Internal and forced climate
 1301 variability during the last millennium: a model-data comparison using ensemble simulations,
 1302 *Quat. Science Rev.*, 24, 1345-1360, 2005.

- 1303 Goosse H., Arzel, O., Luterbacher, J., Mann, M. E., Renssen, H., Riedwyl, N., Timmermann,
1304 A., Xoplaki, E., and Wanner, H.. The origin of the European "Medieval Warm Period". *Clim.*
1305 *Past* 2, 99-113, 2006.
- 1306 Goosse, H., Braidia, M., Crosta, X., Mairesse, A., Masson-Delmotte, V., Mathiot, P., Neukom,
1307 R., Oerter, H., Philippon, G., Renssen, H., Stenni, B., van Ommen, T., and Verleyen, E.:
1308 Antarctic temperature changes during the last millennium: evaluation of simulations and
1309 reconstructions. *Quaternary Science Reviews* 55: 75-90, 2012a.
- 1310 Goosse H., Crespin, E, Dubinkina, S., Loutre, M.F., Mann, M. E, Renssen, H., Sallaz-Damaz,
1311 Y., and Shindell, D.: The role of forcing and internal dynamics in explaining the “Medieval
1312 Climate Anomaly”, *Clim. Dyn.* 39, 2847–2866, doi: 10.1007/s00382-012-1297-0, 2012b.
- 1313 Hansen, J., and Sato, M.: Greenhouse gas growth rates, *Proc. Natl. Acad. Sci.*, 101, 16109-
1314 16114, doi:10.1073/pnas.0406982101, 2004.
- 1315 Hargreaves, J. C., Paul, A., Ohgaito, R., Abe-Ouchi, A., and Annan, J. D.: Are paleoclimate
1316 model ensembles consistent with the MARGO data synthesis? *Climate of the Past* 7, 917–
1317 933, doi:10.5194/cp-7-917-2011, 2011.
- 1318 Hargreaves, J. C., Annan, J. D., Ohgaito, R., Paul, A., and Abe-Ouchi, A.: Skill and reliability
1319 of climate model ensembles at the Last Glacial Maximum and mid-Holocene. *Clim. Past*, 9,
1320 811-823, doi:10.5194/cp-9-811-2013, 2013.
- 1321 Hegerl, G. C., Crowley, T. J., Baum, S. K., Kim, K.-Y., and Hyde, W. T.: Detection of
1322 volcanic, solar and greenhouse gas signals in paleo-reconstructions of Northern Hemispheric
1323 temperature, *Geophys. Res. Lett.*, 30, 1242, doi:10.1029/2002GL016635, 5, 2003.
- 1324 Hegerl, G. C., Crowley, T. J., Hyde, W. T. and Frame, D. J.: Climate sensitivity constrained by
1325 temperature reconstructions over the past seven centuries, *Nature* 440: 1029-1032, 2006.
- 1326 Hegerl, G. C., Crowley, T. J., Allen, M., Hyde, W. T., Pollack, H. N., Smerdon, J., and Zorita,
1327 E.: Detection of human influence on a new, validated 1500-year temperature reconstruction, *J.*
1328 *Clim.*, 20, 650-666, doi:10.1175/JCLI4011.1, 2007.
- 1329 Hegerl, G., Luterbacher, J., González-Rouco, F., Tett, S., Crowley, T., and Xoplaki, E.:
1330 Influence of human and natural forcing on European seasonal temperatures, *Nat. Geoscience*
1331 4 (1179): 99–103, 2011.
- 1332 Hegerl, G.C., and Zwiers, F.W.: Use of models in detection and attribution of climate change,
1333 *WIRES: Climate Change*, 2, 570-591, 2011.
- 1334 Hind, A., and Moberg, A., and Sundberg, R.: Statistical framework forevaluation of climate
1335 model simulations by use of climate proxy data from the last millennium – Part 2: A pseudo-
1336 proxy study addressing the amplitude of solar forcing, *Clim. Past*, 8, 1355– 1365,
1337 doi:10.5194/cp-8-1355-2012, 2012.
- 1338 Hind, A., and Moberg, A. : Past millennial solar forcing magnitude. A statistical hemispheric-
1339 scale climate model versus proxy data comparison, *Clim. Dyn.*, 41, 2527–2537,
1340 doi:10.1007/s00382-012-1526-6, 2013.
- 1341 Hurtt, G. C., Chini, L. P., Frohking, S., Betts, R. A., Feddema, J., Fischer, G., and Wang, Y. P.:
1342 Harmonization of land-use scenarios for the period 1500–2100: 600 years of global gridded

- 1343 annual land-use transitions, wood harvest, and resulting secondary lands, *Clim. Change*,
 1344 109(1-2), 117-161. Doi: 10.1007/s10584-011-0153-2, 2011
- 1345 Jansen, E., Overpeck, J., Briffa, K.R., Duplessy, J.-C., Joos, F., Masson-Delmotte, V., Olago,
 1346 D., Otto-Bliesner, B., Peltier, W.R., Rahmstorf, S., Ramesh, R., Raynaud, D., Rind, D.,
 1347 Solomina, O., Villalba, R., and Zhang, D.: Palaeoclimate. In: *Climate Change 2007: The*
 1348 *Physical Science Basis. Contribution of Working Group I to the Fourth Assessment Report of*
 1349 *the Intergovernmental Panel on Climate Change [Solomon, S., D. Qin, M. Manning, Z. Chen,*
 1350 *M. Marquis, K.B. Averyt, M. Tignor and H.L. Miller (eds.)]. Cambridge University Press,*
 1351 *Cambridge, United Kingdom and New York, NY, USA, 2007.*
- 1352 Johns, T. C., Gregory, J. M., Ingram, W. J., Johnson, C. E., Jones, A., Lowe, J. A., and
 1353 Woodage, M. J.: Anthropogenic climate change for 1860 to 2100 simulated with the HadCM3
 1354 model under updated emissions scenarios, *Clim. Dyn.*, 20(6), 583-612, doi: 10.1007/s00382-
 1355 002-0296-y, 2003
- 1356 Jones, P.D., Briffa, K.R., Osborn, T.J., Lough, J. M., van Ommen, T., Vinther, B.M.,
 1357 Luterbacher, J., Zwiers, F.W., Wahl, E., Schmidt, G., Ammann, C., Mann, M.E., Wanner, H.,
 1358 Buckley, B.M., Cobb, K., Esper, J., Goosse, H., Graham, N., Jansen, E., Kiefer, T., Kull, C.,
 1359 Mosley-Thompson, E., Overpeck, J.T., Schulz, M., Tudhope, S., Villalba, R., and Wolff, E.:
 1360 High-resolution paleoclimatology of the last millennium: a review of the current status and
 1361 future prospects, *The Holocene* 19, 3-49, 2009.
- 1362 Jones, P.D. and Mann, M.E.: Climate over past millennia, *Rev. Geophys.* 42, RG2002, doi:
 1363 10.1029/2003RG000143, 2004.
- 1364 Jungclauss, J. H., Keenlyside, N., Botzet, M., Haak, H., Luo, J.-J., Latif, M., Marotzke, J.,
 1365 Mikolajewicz, U., and Roeckner, E.: Ocean Circulation and Tropical Variability in the
 1366 Coupled Model ECHAM5/MPI-OM, *J. Clim.*, 19, 3952– 3972, 2006.
- 1367 Jungclauss, J. H., Lorenz, S. J., Timmreck, C., Reick, C. H., Brovkin, V., Six, K., and
 1368 Marotzke, J.: Climate and carbon-cycle variability over the last millennium. *Climate of the*
 1369 *Past*, 6, 723-737. doi: 10.5194/cp-6-723-2010, 2010.
- 1370 Jungclauss, J. H., Lohmann, K., and Zanchettin, D. Enhanced 20th-century heat transfer to the
 1371 Arctic simulated in the context of climate variations over the last millennium. *Clim. Past*,
 1372 10(6), 2201-2213, doi: 10.5194/cp-10-2201-2014, 2014.
- 1373 Kaplan, J. O., Krumhardt, K. M., Ellis, E. C., Ruddiman, W. F., Lemmen, C., & Goldewijk, K.
 1374 K.: Holocene carbon emissions as a result of anthropogenic land cover change. *The Holocene*,
 1375 21(5) 775,–79, doi: 10.1177/0959683610386983, 2011.
- 1376 Knight, J.R., Allan R.J., Folland C.K., Vellinga M., and Mann M.E.: A signature of persistent
 1377 natural thermohaline circulation cycles in observed climate. *Geophys. Res. Let.* 32 (20),
 1378 L20708, 2005.
- 1379 Krivova, N. A., Balmaceda, L., and Solanki, S. K.: Reconstruction of solar total irradiance
 1380 since 1700 from the surface magnetic flux, *Astron. Astrophys.* 467(1), 335-346, doi:
 1381 10.1051/0004-6361:20066725, 2007
- 1382 Kutzbach, L., Thees, B., and Wilmking, M.: Identification of linear relationships from noisy
 1383 data using errors-in-variables models-relevance for reconstruction of past climate from tree-

- 1384 ring and other proxy information, *Clim. Change*, 105, 155–177, doi 10.1007/s10584-010-
1385 9877-7, 2011.
- 1386 Laepple, T., and Huybers, P.: Global and regional variability in marine surface temperatures,
1387 *Geophys. Res. Lett.*, 41, doi:10.1002/2014GL059345, 2014a.
- 1388 Laepple, T., and Huybers, P.: Ocean surface temperature variability: Large model–data
1389 differences at decadal and longer periods, *Proceed. Nat. Acad. Sciences* 11 (47) 16682–16687,
1390 doi: 10.1073/pnas.1412077111, 2014b.
- 1391 Lamarque, J. F., Bond, T. C., Eyring, V., Granier, C., Heil, A., Klimont, Z., and Van Vuuren,
1392 D. P.: Historical (1850–2000) gridded anthropogenic and biomass burning emissions of
1393 reactive gases and aerosols: methodology and application., *Atmos. Chem. Phys.*, 10(15),
1394 7017-7039 doi: 10.5194/acp-10-7017-2010, 2010
- 1395 Landrum, L., Otto-Bliesner, B. L., Wahl, E. R., Conley, A., Lawrence, P. J., Rosenbloom, N.,
1396 and Teng, H.: Last millennium climate and its variability in CCSM4, *J. Clim.* 26(4), 1085-
1397 1111, doi: 10.1175/JCLI-D-11-00326.1, 2013
- 1398 Lehner, F., Raible, C. C., and Stocker, T. F.: Testing the robustness of a precipitation proxy-
1399 based North Atlantic Oscillation reconstruction, *Quat. Sci. Rev.*, 45, 85-94, 2012.
- 1400 Lehner, F., Born, A., Raible, C. C., Stocker, T. F.: Amplified inception of European Little Ice
1401 Age by sea ice-ocean-atmosphere feedbacks, *J. Clim.* 26, 7586-7602, doi: 10.1175/JCLI-D-
1402 12-00690.1, 2013.
- 1403 Lehner, F., Joos, F., Raible, C. C., Mignot, J., Born, A., Keller, K. M., and Stocker, T. F.:
1404 Climate and carbon cycle dynamics in a CESM simulation from 850-2100 CE, *Earth Sys.*
1405 *Dyn. Discuss.*, 6, 1–56, 2015.
- 1406 Lohmann, K., Jungclauss, J. H., Matei, D., Mignot, J., Menary, M., Langehaug, H. R., Ba, J.,
1407 Gao, Y., Otterå, O. H., Park, W., and Lorenz, S.: The role of subpolar deep water formation
1408 and Nordic Seas overflows in simulated multidecadal variability of the Atlantic meridional
1409 overturning circulation, *Ocean Sci.*, 10, 227-241, doi:10.5194/os-10-227-2014, 2014
- 1410 Lorenz, E.: Deterministic Nonperiodic Flow, *J. Atmos. Sci.*, 20, 130–141, 1963.
- 1411 Luterbacher, J., Dietrich, D., Xoplaki, E., Grosjean, M., and Wanner, H.: European seasonal
1412 and annual temperature variability, trends, and extremes since 1500, *Science* 5: 303, 1499-
1413 1503, doi:10.1126/science.1093877, 2004.
- 1414 MacFarling Meure, C., Etheridge, D., Trudinger, C., Steele, P., Langenfelds, R., Van Ommen,
1415 T., and Elkins, J.: Law Dome CO₂, CH₄ and N₂O ice core records extended to 2000 years
1416 BP. *Geophysical Research Letters*, 33(14). Doi: 10.1029/2006GL026152, 2006.
- 1417 Machida, T., Nakazawa, T., Fujii, Y., Aoki, S., and Watanabe, O.: Increase in the atmospheric
1418 nitrous oxide concentration during the last 250 years, *Geophys. Res. Lett.*, 22(21), 2921-2924.
1419 Doi: 10.1029/2001GB001417, 1995.
- 1420 Mann, M.E., Bradley, R.S., and Hughes, M.K.: Northern Hemisphere temperatures during the
1421 past millennium: inferences, uncertainties, and limitations, *Geophys. Res. Lett.* 26, 759-762,
1422 1999.

- 1423 Mann, M. E., Fuentes, J. D. and Rutherford, S.: Underestimation of volcanic cooling in tree-
1424 ring-based reconstructions of hemispheric temperatures, *Nat. Geosci.*, 5(3), 202–205, 2012.
- 1425 Mann, M., Zhang, Z., Rutherford, S., Bradley, R., Hughes, M., Shindell, D., Ammann, C.,
1426 Faluvegi, G. and Ni, F.: Global signatures and dynamical origins of the Little Ice Age and
1427 Medieval Climate Anomaly, *Science* 326, 1256–1260, 2009.
- 1428 Marzban, C., Wang, R., Kong, F., and Leyton, S.: On the effect of correlations on rank
1429 histograms: Reliability of temperature and wind speed forecasts from finescale ensemble
1430 reforecasts, *Month. Weath. Rev.*, 139(1), 295-310, 2011.
- 1431 Masson-Delmotte, V., Schulz, M., Abe-Ouchi, A., Beer, J., Ganopolski, A., González Rouco,
1432 J. F., Jansen, E., Lambeck, K., Luterbacher, J., Naish, T., Osborn, T., Otto-Bliesner, B., Quinn,
1433 T., Ramesh, R., Rojas, M., Shao, X. and Timmermann, A.: Information from Paleoclimate
1434 Archives. In: *Climate Change 2013: The Physical Science Basis. Contribution of Working
1435 Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*
1436 [Stocker, T. F., D. Qin, G.-K. Plattner, M. Tignor, S. K. Allen, J. Boschung, A. Nauels, Y. Xia,
1437 V. Bex and P. M. Midgley (eds.)]. Cambridge University Press, Cambridge, United Kingdom
1438 and New York, NY, USA, 2013.
- 1439 McGregor H.V., Evans, M. N., Goosse, H., Leduc, G., Martrat, B., Addison, J. A., Mortyn, P.
1440 G., Oppo, D. W., Seidenkrantz, M.S., Sicre, M.-A., Phipps, S. J., Selvaraj, K., Thirumalai, K.,
1441 Filipsson, H. L., and Ersek, V.: Robust global ocean cooling trend for the past two millennia.
1442 *Nat. Geos.*, 8(9), 671-677, DOI:10.1038/NNGEO2510, 2015.
- 1443 McKay, N.P., and Kaufman, D.S.: An extended Arctic proxy temperature database for the past
1444 2,000 years. *Scientific Data* 1:140026 doi: 10.1038/sdata.2014.26, 2014.
- 1445 Miller, G., Geirsdóttir, Á., Zhong, Y., Larsen, D. J., Otto-Bliesner, B. L., Holland, M. M.,
1446 Bailey, D. A., Refsnider, K. A., Lehman, S. J., Southon, J. R., Anderson, C., Björnsson, H., and
1447 Thordarson, T.: Abrupt onset of the Little Ice Age triggered by volcanism and sustained by
1448 sea-ice/ocean feedbacks, *Geophys. Res. Lett.*, 39, L02708, doi:10.1029/2011GL050168, 2012.
- 1449 Moberg, A., Mohammad, R., and Mauritsen, T.: Analysis of the Moberg et al. (2005)
1450 hemispheric temperature reconstruction, *Clim. Dyn.* 31, 957-971, doi:10.1007/s00382-008-
1451 0392-8, 2008
- 1452 Moberg, A., and Brattström, G.: Prediction intervals for climate reconstructions with
1453 autocorrelated noise—An analysis of ordinary least squares and measurement error methods,
1454 *Palaeoge., Palaeoclim., Palaeoecol.* 308, 313–329, 2011
- 1455 Moberg, A.: Comparisons of simulated and observed Northern Hemisphere temperature
1456 variations during the past millennium – selected lessons learned and problems encountered,
1457 *Tellus B* 65, 19921, doi:10.3402/tellusb.v65i0.19921, 2013.
- 1458 Moberg, A., Sundberg, R., Grudd H., and Hind; A.: Statistical framework for evaluation of
1459 climate model simulations by use of climate proxy data from the last millennium – Part 3:
1460 Practical considerations, relaxed assumptions, and using tree-ring data to address the
1461 amplitude of solar forcing, *Clim. Past* 11, 425–448, doi:10.5194/cp-11-425-2015, 2015

- 1462 Morice, C.P., Kennedy, J.J., Rayner, N.A. and Jones, P.D.: Quantifying uncertainties in global
1463 and regional temperature change using an ensemble of observational estimates: the
1464 HadCRUT4 dataset, *J. Geophys. Res.*, 117, D08101, doi:10.1029/2011JD017187, 2012.
- 1465 Murphy, A. H.: A new vector partition of the probability score, *J. Appl. Meteorol.*, 12, 595–
1466 600, 1973.
- 1467 Neukom, R., and Gergis, J.: Southern Hemisphere high-resolution palaeoclimate records of
1468 the last 2000 years, *The Holocene* 5: 501–524, 2012.
- 1469 Neukom, R., Gergis, J., Karoly, D., Wanner, H., Curran, M., Elbert, J., González-Rouco, F.,
1470 Linsley, B., Moy, A., Mundo, I., Raible, C., Steig, E., van Ommen, T., Vance, T., Villalba, R.,
1471 Zinke, J. and Frank, D.: Inter-hemispheric temperature variability over the last millennium,
1472 *Nat. Clim. Change* 4: 362–367, 2014
- 1473 Ortega, P., Lehner, F., Casado, M., Swingedouw, D., Masson-Delmotte, V., Yiou, P., Raible, C.
1474 C.: A multi-proxy model-tested NAO reconstruction for the last millennium. *Nature*,
1475 submitted, 2015.
- 1476 PAGES 2k Consortium : Ahmed, M. , Anchukaitis, K., Asrat, A., Borgaonkar, H., Braidia, M.,
1477 Buckley, B. , Büntgen, U., Chase, B., Christie, D., Cook, E., Curran, M., Diaz, H. , Esper, J.,
1478 Fan, Z.X., Gaire, N., Ge, Q., Gergis, J. , Gonzalez-Rouco, J.F., Goosse, H. , Grab, S., Graham,
1479 R., Graham, N., Grosjean, M., Hanhijärvi, S., Kaufman, D., Kiefer, T., Kimura, K., Korhola,
1480 A., Krusic, P., Lara, A., Lézine, A.M., Ljungqvist, F., Lorrey, A., Luterbacher, J., Masson-
1481 Delmotte, D. McCarroll, J. McConnell, N. McKay, M. Morales, A. Moy, R. Mulvaney, I.
1482 Mundo, V., Nakatsuka, T., Nash, D., Neukom, R., Nicholson, S., Oerter, H., Palmer, J.,
1483 Phipps, S., Prieto, M., Rivera, A., Sano, M., Severi, M., Shanahan, T., Shao, X., Shi, F., Sigl,
1484 M., Smerdon, J., Solomina, O., Steig, E., Stenni, B., Thamban, M., Trouet, V., Turney, C.,
1485 Umer, M., van Ommen, T., Verschuren, D., Viau, A., Villalba, R., Vinther, B., von Gunten, L.,
1486 Wagner, S., Wahl, E., Wanner, H., Werner, J., White, J., Yasue, K., Zorita, E.: Continental-
1487 scale temperature variability during the last two millennia. *Nature Geoscience* 6, 339-346
1488 DOI: 10.1038/NGEO1797, 2013.
- 1489 PAGES 2k Consortium (Primary Authors: K. Anchukaitis, U. Büntgen, J. Emile-Geay, M. N.
1490 Evans, H. Goosse, D. Kaufman, J. Luterbacher, J. Smerdon, M. Tingley, L. von Gunten): A
1491 Community-driven framework for climate reconstructions, *Eos Trans. AGU*, 95(40), 361,
1492 2014.
- 1493 Percival, D.B., and Walden, A.T., *Spectral Analysis for Physical Applications: Multitaper and*
1494 *Conventional Univariate Techniques*, Cambridge University Press, 1993.
- 1495 Phipps, S. J., McGregor, H. V., Gergis, J., Gallant, A. J., Neukom, R., Stevenson, S. and Van
1496 Ommen, T. D.: Paleoclimate data–model comparison and the role of climate forcings over the
1497 past 1500 years, *J. Clim.*, 26(18), 6915-6936, doi: 10.1175/JCLI-D-12-00108.1, 2013.
- 1498 Pongratz, J., Reick, C., Raddatz, T., and Claussen, M.: A reconstruction of global agricultural
1499 areas and land cover for the last millennium, *Glob. Biogeochem. Cycles*, 22(3), doi:
1500 10.1029/2007GB003153, 2008.
- 1501 Pongratz, J., Reick, C. H., Raddatz, T., and Claussen, M.: Effects of anthropogenic land cover
1502 change on the carbon cycle of the last millennium, *Glob. Biogeochem. Cycles*, 23(4), doi:
1503 10.1029/2009GB003488, 2009.

- 1504 Raible, C. C., Stocker, T. F., Yoshimori, M., Renold, M., Beyerle, U., Casty, C., and
 1505 Luterbacher, J.: Northern Hemispheric trends of pressure indices and atmospheric circulation
 1506 patterns in observations, reconstructions, and coupled GCM simulations, *J. Climate*, 18, 3968-
 1507 3982, 2005.
- 1508 Raible, C.C., C. Casty, J. Luterbacher, A. Pauling, J. Esper, D.C. Frank, U. Büntgen, A.C.
 1509 Roesch, P. Tschuck, M. Wild, P-L. Vidale, C. Schär, and H. Wanner: Climate Variability -
 1510 Observations, Reconstructions, and Model Simulations for the Atlantic-European and Alpine
 1511 region from 1500-2100 AD, *Clim. Change*, 79, 9-29, 2006.
- 1512 Raible, C. C., Lehner, F., González-Rouco, J. F., and Fernández-Donado, L.: Changing
 1513 correlation structures of the Northern Hemisphere atmospheric circulation from 1000 to 2100
 1514 AD, *Clim. Past*, 10, 537-550, doi:10.5194/cp-10-537-2014, 2014.
- 1515 Rougier, J., Goldstein, M., and House, L.: Second-order exchangeability analysis for
 1516 multimodel ensembles. *J. Am. Stat. Assoc.*, 108(503), 852-863, 2013.
- 1517 Russon, T., Tudhope A., Hegerl, G.C., and Collins M.: Inferring changes in ENSO amplitude
 1518 from proxy records, *Geophys. Res. Lett.*, in press, 2015.
- 1519 Schmidt, G. A., Jungclaus, J. H., Ammann, C. M., Bard, E., Braconnot, P., Crowley, T. J.,
 1520 Delaygue, G., Joos, F., Krivova, N. A., Muscheler, R., Otto-Bliesner, B. L., Pongratz, J.,
 1521 Shindell, D. T., Solanki, S. K., Steinhilber, F., and Vieira, L. E. A.: Climate forcing
 1522 reconstructions for use in PMIP simulations of the last millennium (v1.0), *Geosci. Model*
 1523 *Dev.*, 4, 33-45, doi:10.5194/gmd-4-33-2011, 2011.
- 1524 Schmidt, G. A., Jungclaus, J. H., Ammann, C. M., Bard, E., Braconnot, P., Crowley, T. J.,
 1525 Delaygue, G., Joos, F., Krivova, N. A., Muscheler, R., Otto-Bliesner, B. L., Pongratz, J.,
 1526 Shindell, D. T., Solanki, S. K., Steinhilber, F., and Vieira, L. E. A.: Climate forcing
 1527 reconstructions for use in PMIP simulations of the last millennium (v1.1), *Geosci. Model*
 1528 *Dev.*, 5, 185–191, doi:10.5194/gmd-5-185-2012, 2012.
- 1529 Schmidt, G. A., Annan, J. D., Bartlein, P. J., Cook, B. I., Guilyardi, E., Hargreaves, J. C.,
 1530 Harrison, S. P., Kageyama, M., LeGrande, A. N., Konecky, B., Lovejoy, S., Mann, M. E.,
 1531 Masson-Delmotte, V., Risi, C., Thompson, D., Timmermann, A., Tremblay, L.-B., and
 1532 Yiou, P.: Using palaeo-climate comparisons to constrain future projections in CMIP5, *Clim.*
 1533 *Past*, 10, 221-250, doi:10.5194/cp-10-221-2014, 2014a.
- 1534 Schmidt, G. A., Kelley, M., Nazarenko, L., Ruedy, R., Russell, G. L., Aleinov, I., and Zhang,
 1535 J.: Configuration and assessment of the GISS ModelE2 contributions to the CMIP5 archive.
 1536 *Journal of Advances in Modeling Earth Systems*, 6(1), 141-184. doi:
 1537 10.1002/2013MS000265, 2014b.
- 1538 Schurer A.P., Hegerl G.C., Mann M. E., Tett, S. F. B., and Phipps, S. J.: Separating forced
 1539 from chaotic climate variability over the Past Millennium, *J. Climate*, 26, 6954–6973, 2013.
- 1540 Schurer, A. P., Tett, S. F. and Hegerl, G. C.: Small influence of solar variability on climate
 1541 over the past millennium. *Nat. Geoscience*, 7 (2), 104–108, doi:10.1038/NGEO2040, 2014.
- 1542 Shapiro, A. I., Schmutz, W., Rozanov, E., Schoell, M., Haberreiter, M., Shapiro, A. V., and
 1543 Nyeki, S.: A new approach to the long-term reconstruction of the solar irradiance leads to

- 1544 large historical solar forcing, *Astron. Astrophys.*, 529, A67, doi:10.1051/0004-
1545 6361/201016173, 2011
- 1546 Shi, F., Yang, B., Mairesse, A., von Gunten, L., Li, J., Bräuning, A., Yang, F., Xiao, X.:
1547 Northern Hemisphere temperature reconstruction during the last millennium using multiple
1548 annual proxies. *Climate Research*, 56(3), 231-244, 2013.
- 1549 Shi, F., Ge, Q., Yang, B., Li, J., Yang, F., Charpentier Ljungqvist, F., Solomina, O., Nakatsuka,
1550 T., Wang, N., Zhao, S., Xu, C., Fang, K., Sano, M., Chu, G., Fan, Z., Gaire, N. P., Zafar, M.
1551 U.: A multi-proxy reconstruction of spatial and temporal variations in Asian summer
1552 temperatures over the last millennium. *Climatic Change*, 131(4): 663-676, 2015.
- 1553 Shindell, D.T., Schmidt, G.A., Mann, M.E., Rind, D., and Waple, A., Solar forcing of regional
1554 climate change during the Maunder Minimum, *Science*, 294, 2149-2152, 2001.
- 1555 Smerdon, J. E., Kaplan, A., Chang, D., and Evans, M. N.: A pseudoproxy evaluation of the
1556 CCA and RegEM methods for reconstructing climate fields of the last millennium, *J. Clim.*,
1557 23, 4856–4880, 2010.
- 1558 Smerdon, J. E.: Climate models as a test bed for climate reconstruction methods: pseudoproxy
1559 experiments, *WIRES Climate Change*, 3, 63-77 doi: 10.1002/wcc.149, 2012.
- 1560 Smerdon, J.E, Cook, B.I., Cook, E.R., and Seager, R.: Bridging past and future climate across
1561 paleoclimatic reconstructions, observations, and models: a hydroclimate case study, *J. Clim.*,
1562 28(8), 3212-3231, 2015a.
- 1563 Smerdon, J.E., S. Coats, and T.R. Ault, Model-Dependent Spatial Skill in Pseudoproxy
1564 Experiments Testing Climate Field Reconstruction Methods for the Common Era, *Climate*
1565 *Dynamics*, (submitted), 2015b.
- 1566 Steiger, N. J., Hakim, G., Steig, E. J., Battisti, D. S. and Roe, G. H.: Assimilation of time-
1567 averaged pseudoproxies for climate, *J. Clim.*, 27, 426-441, 2014.
- 1568 Steinhilber, F., Beer, J., and Fröhlich, C.: Total solar irradiance during the Holocene.
1569 *Geophysical Research Letters*, 36(19), doi: 10.1029/2009GL040142, 2009
- 1570 Stenchikov, G., Hamilton, K., Stouffer, R. J., Robock, A., Ramaswamy, V., Santer, B., and
1571 Graf, H. F.: Arctic Oscillation response to volcanic eruptions in the IPCC AR4 climate
1572 models, *J. Geophys. Res.-Atmos.*, 111, D07107, doi:10.1029/2005JD006286, 2006.
- 1573 Stothers, R.B.: The great Tambora eruptions in 1815 and its aftermath. *Science*, 224(4654),
1574 1191-1198, 1984.
- 1575 Sundberg, R., Moberg, A., Hind A.: Statistical framework for evaluation of climate model
1576 simulations by use of climate proxy data from the last millennium – Part 1: Theory, *Clim. Past*
1577 8, 1339-1353, doi:10.5194/cp-8-1339-2012, 2012.
- 1578 Taylor, K. E., Stouffer, R. J., and Meehl, G. A.: An overview of CMIP5 and the experiment
1579 design,, *Bull. Amer. Meteor. Soc.*, 93, doi:10.1175/BAMS-D-11-00094.1., 2012.
- 1580 Tett, S. F. B., Betts, R., Crowley, T. J., Gregory, J., Johns, T. C., Jones, A., Osborn, T. J.,
1581 Ostrom, E., Roberts, D. L., and Woodage, M. J.: The impact of natural and anthropogenic
1582 forcings on climate and hydrology since 1550, *Clim. Dyn.*, 28(1), 3–34, 2007

1583 Tierney, J. E., Abram, N. J., Anchukaitis, K. J., Evans, M. N., Giry, C., Kilbourne, K. H.,
1584 Saenger, C. P., Wu, H. C., and Zinke, J.: Tropical sea-surface temperatures for the past four
1585 centuries reconstructed from coral archives, *Paleoceanography*,30, 226–252,
1586 doi:10.1002/2014PA002717, 2015.

1587 Tingley, M. P., Craigmile, P. F., Haran, M., Li, B., Mannshardt, E., and Rajaratnam, B.:
1588 Piecing together the past: statistical insights into paleoclimatic reconstructions, *Quat. Science*
1589 *Rev.*, 35, 1-22, 2012.

1590 Tingley, M. P., Stine, A. R., and Huybers, P. : Temperature reconstructions from tree-ring
1591 densities overestimate volcanic cooling, *Geophys. Res. Lett.*, 41, 7838–7845,
1592 doi:10.1002/2014GL061268, 2014.

1593 Thomson, D.J.: Spectrum estimation and harmonic analysis, *Proceed. IEEE.*, 70,1055-1096,
1594 1982.

1595 Vieira, L. E. A., Solanki, S. K., Krivova, N. A., and Usoskin, I.: Evolution of the solar
1596 irradiance during the Holocene, *Astron. Astrophys.*, 531, A6, doi: 10.1051/0004-
1597 6361/201015843, 2011.

1598 Wang, J., Emile-Geay, J., Guillot, D., Smerdon, J.E., and Rajaratnam, B.: Evaluating climate
1599 field reconstruction techniques using improved emulations of real-world conditions, *Clim.*
1600 *Past*, 10, 1-19, doi:10.5194/cp-10-1-2014, 2014.

1601 Wang, J., Emile-Geay, J., Guillot, D., McKay, N. P., and B. Rajaratnam (2015), Fragility of
1602 reconstructed temperature patterns over the Common Era: Implications for model evaluation,
1603 *Geophys. Res. Lett.*, 42, 7162–7170, doi:10.1002/2015GL065265.

1604 Wang, Y. M., Lean, J. L., and Sheeley Jr, N. R.: Modeling the sun's magnetic field and
1605 irradiance since 1713, *Astrophys. J.*, 625(1), 522. Doi: 10.1086/429689, 2005.

1606 Widmann, M., Goosse, H., van der Schrier, G., Schnur, R., and Barkmeijer, J.: Using data
1607 assimilation to study extratropical Northern Hemisphere climate over the last millennium.
1608 *Climate of the Past* 6, 627–644, 2010

1609 Wilmes, S., Raible, C. and Stocker, T.: Climate variability of the mid- and high-latitudes of
1610 the Southern Hemisphere in ensemble simulations from 1500 to 2000AD, *Clim. Past* 8: 373–
1611 390, 2012.

1612 Wigley, T. M. L., Ammann, C. M., Santer, B. D., and Raper, S. C. B.: Effect of climate
1613 sensitivity on the response to volcanic forcing, *J. Geophys. Res.*, 110, D09107,
1614 doi:10.1029/2004JD005557, 2005.

1615 Wunsch, C.: The interpretation of short climate records, with comments on the North Atlantic
1616 and Southern Oscillations, *Bull. Amer. Meteor. Soc.*, 80, 245-255, 1999.

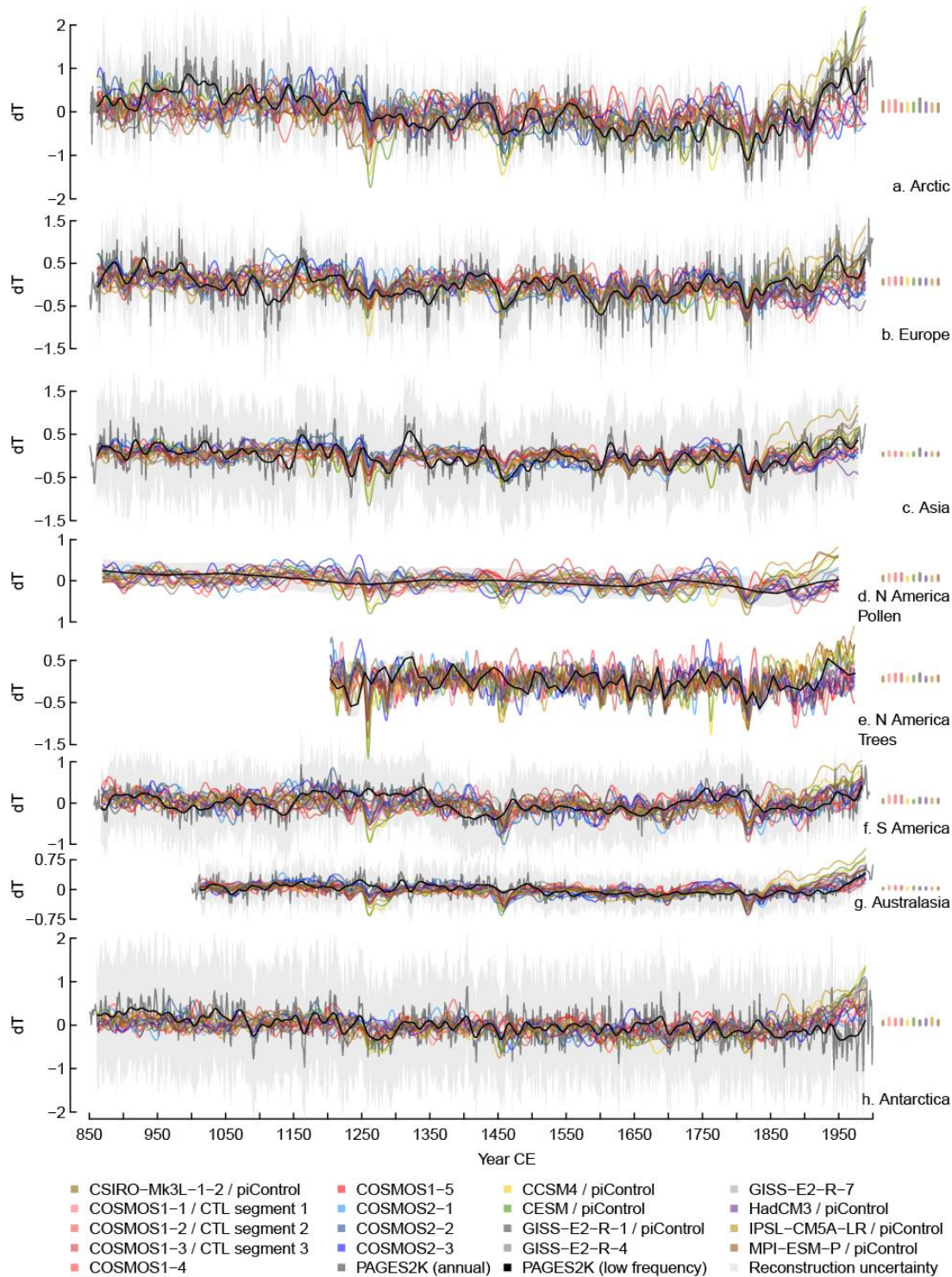
1617 Yoshimori, M., Stocker, T. F., Raible, C. C., and Renold, M.: Externally-forced and internal
1618 variability in ensemble climate simulations of the Maunder Minimum, *J. Climate*, 18, 4253-
1619 4270, 2005.

1620 Zanchettin, D., Bothe, O., Lehner, F., Ortega, P., Raible, C. C., and Swingedouw, D.:
1621 Reconciling reconstructed and simulated features of the winter Pacific-North-American
1622 pattern in the early 19th century, *Clim. Past*, submitted, 2015.

1623 Zorita, E., Gonzalez-Rouco, F., and Legutke, S.: Testing the Mann et al. (1998) approach to
1624 paleoclimate reconstructions in the context of a 1000-yr control simulation with the ECHO-G
1625 coupled climate model, *J. Climate*, 16, 1378-1390, 2003.

1626 Zunz V., Goosse, H. and Massonnet, F.: How does internal variability influence the ability of
1627 CMIP5 models to reproduce the recent trend in Southern Ocean sea ice extent? *The*
1628 *Cryosphere* 7, 451–468, 2013.

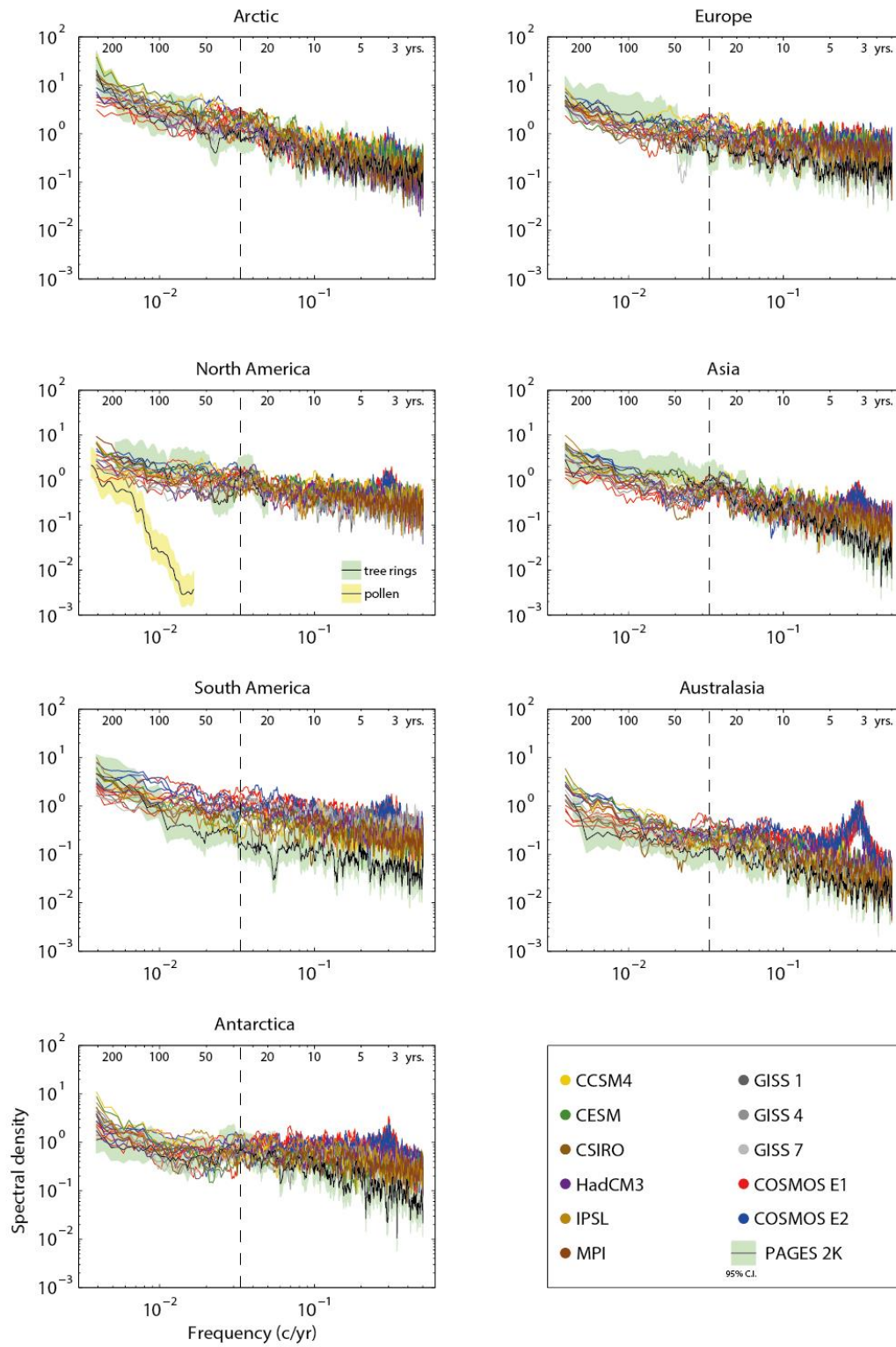
1629



1630

1631 Figure 1: Series of simulated temperatures and PAGES 2k reconstructions for the seven
 1632 continent-scale regions. The reconstructions are shown at their original resolution and after a
 1633 smoothing using a 23-year Hamming filter, except for the North American reconstructions.
 1634 Only the smoothed series are shown for models. Grey shading denotes each reconstruction's
 1635 original uncertainty estimates. Segments on the right indicate the unforced variability of the
 1636 23-year Hamming filtered times series in the respective control simulations (standard
 1637 deviation of the time series, colours as in the caption). The anomalies are computed compared
 1638 to the mean of the time series over the full length of temporal overlap between simulations

1639 and reconstruction. Note the different scales in the y-axis of the various regions.

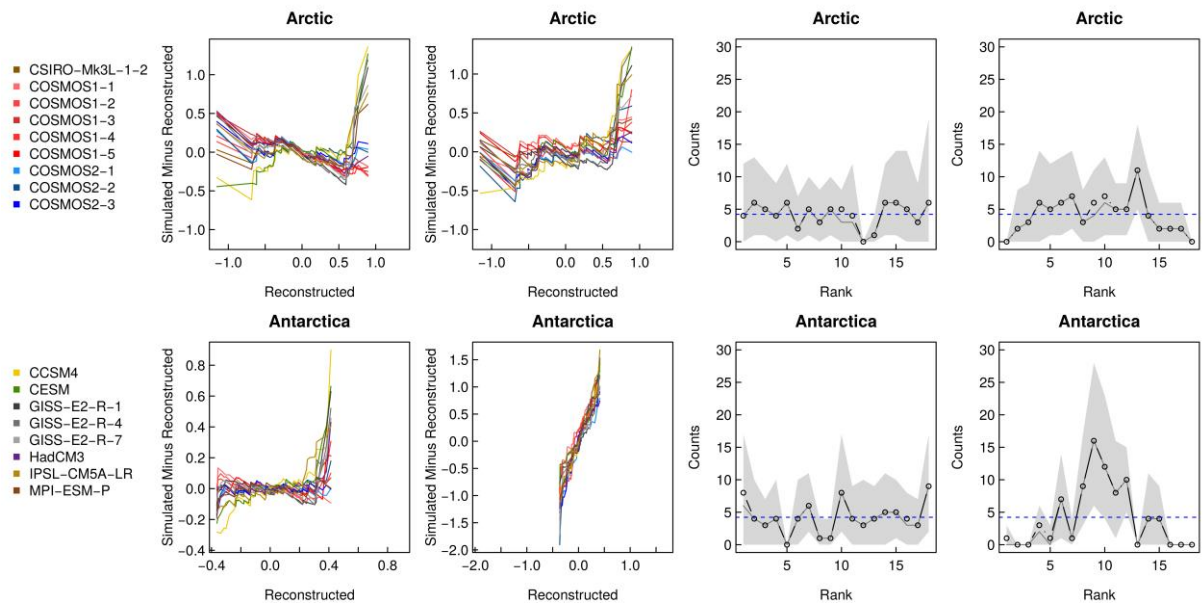


1642

1643 Figure 2: Spectral densities for simulations and reconstructions for PAGES 2k regions,
1644 calculated using all existing data in the period 850-2005 CE. Reconstruction spectra are
1645 illustrated with their 95% confidence intervals in coloured bands, while model spectra are
1646 shown with single coloured lines. Dashed vertical lines denote the limit for frequencies and
1647 periods of relevance (to the left of the line) for analyses made at the 15-year resolution, or
1648 with a 23-point Hamming window, as in many other analyses in this study. The multi-taper
1649 method (Thomson, 1982; Percival and Walden, 1993) was used, with the time-bandwidth
1650 product set to 4 and with long-term averages subtracted before estimating the spectra. Units
1651 are temperature variance ($^{\circ}\text{C}^2$ or K^2) per frequency (c/year).

1652

1653

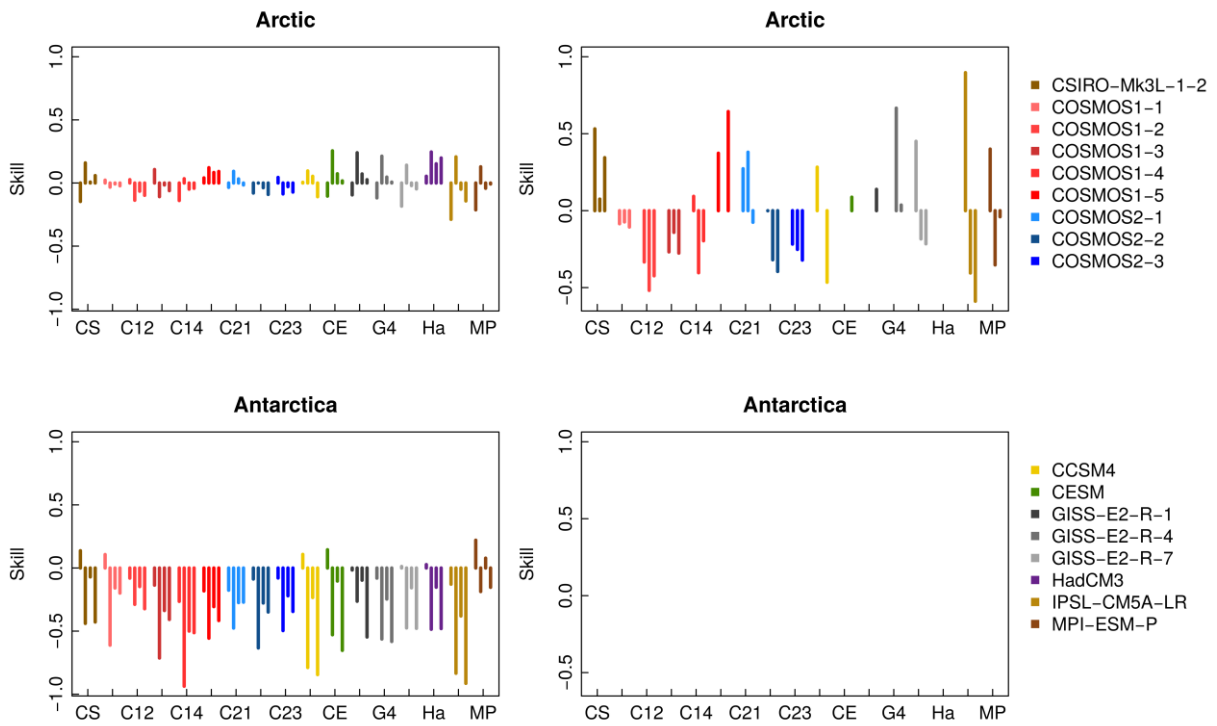


1654

1655 Figure 3: Climatological consistency (first two columns): residual quantile-quantile plots for
1656 the full period; and probabilistic consistency (last 2 columns): rank counts for the full period.
1657 The top row is for the Arctic, and the bottom row is for Antarctica. For both the climatological
1658 and probabilistic consistency, the computations are obtained by neglecting the uncertainties
1659 (left plot) and using the uncertainties provided with the original reconstructions (right plot).
1660 For the climatological assessment, positive and negative slopes or large differences from 0
1661 emphasize lack of consistency. For the probabilistic measure, U- or dome-shaped features
1662 highlight lack of consistency.

1663

1664

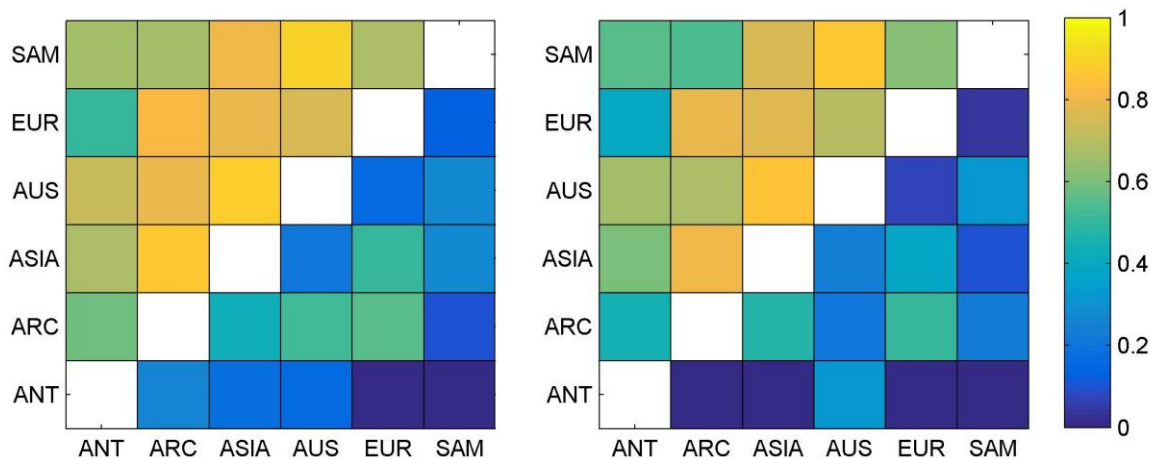


1666

1667 Figure 4: Skill metric for the individual models for all periods (from left to right: 850-1350,
 1668 1350-1850, 850-1850, 850-2000). Top row for the Arctic, bottom for Antarctica. The
 1669 computations assume no uncertainties (left plot) and uncertainties provided with the original
 1670 reconstructions (right plot). When the skill is undefined (as for Antarctica when using the
 1671 original error estimates) no bar is shown. Positive values indicate skill in this simple
 1672 evaluation.

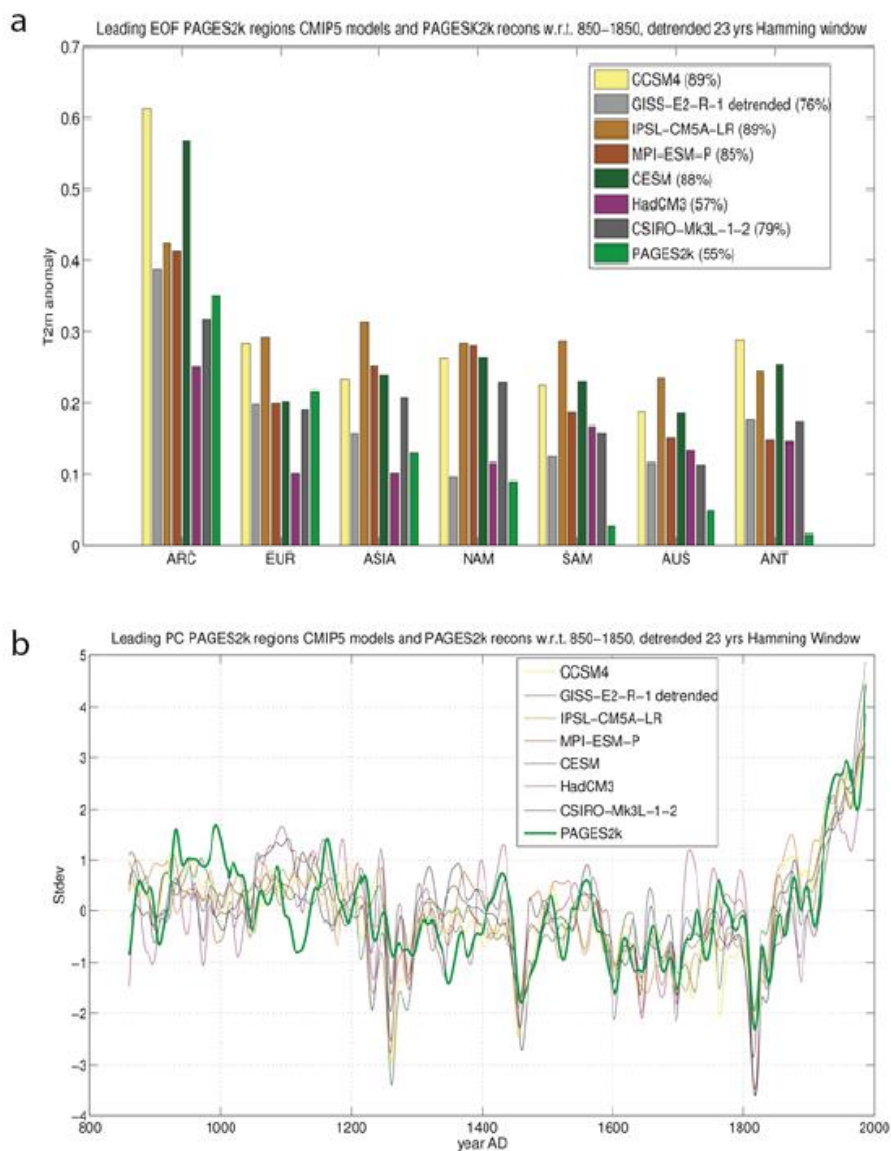
1673

1674



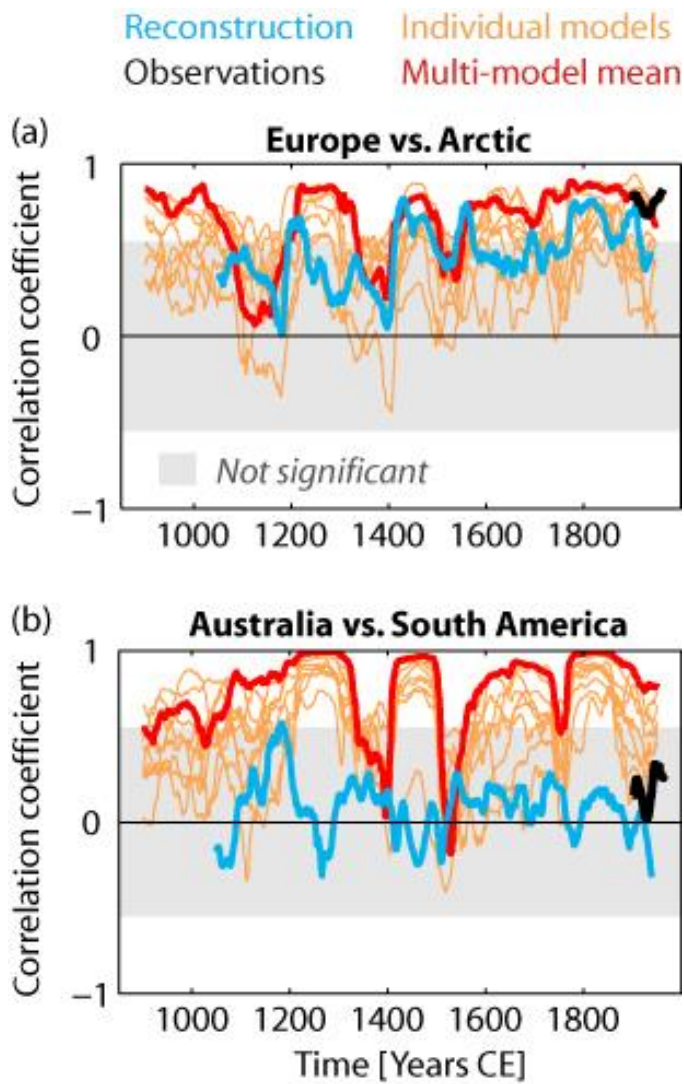
1676
 1677
 1678
 1679
 1680
 1681
 1682
 1683
 1684

Figure. 5: Correlations among the PAGES2k regions for detrended simulated and reconstructed time series filtered using a 23-year Hamming filter. Left-hand panel: forced simulation with MPI-ESM (upper triangle) PAGES 2k reconstructions (lower triangle) for 1012-1978 CE. Right-hand panel: forced simulation with MPI-ESM for the preindustrial period 1012 CE – 1850 CE (upper triangle) and unforced control simulation with MPI-ESM (lower triangle).



1686
 1687 Figure 6: a) Leading EOFs of the near-surface temperature simulated by each CMIP5/PMIP3
 1688 model and reconstructions over the full period 850–2004 CE. The EOF analysis is based on
 1689 the covariance matrix with respect to temperature anomalies for the pre-industrial period 850–
 1690 1850 CE. Values in parentheses correspond to the amount of variance represented by the
 1691 leading EOF. b) Time series of the principal components (PCs) corresponding to the leading
 1692 EOF for the PMIP3 simulations and PAGE2k reconstructions. The time series were filtered
 1693 with a 23-year Hamming filter and were linearly detrended before the covariance matrix was
 1694 calculated. The PC time series are shown as standardized anomalies from the average over the
 1695 full period 850–2004 CE. Positive PC values correspond to positive temperature anomalies in
 1696 the respective regions. Results for single member realizations and the pre-industrial period are
 1697 presented in the Figures S7 and S8, respectively.
 1698
 1699

1700



1701

1702

1703 Figure 7: 100-year moving Tukey window correlations between selected PAGES 2k regions
1704 for the PAGES 2k reconstructions (blue) and PMIP3 models (8 models in orange,
1705 multi model mean in red) and observations from HadCRUT4 (Morice et al., 2012, black).
1706 Each 100-year segment is linearly detrended beforehand. Grey shading illustrates correlations
1707 that are not significant at the 5% level. (a) Correlation between Arctic and Europe as an
1708 example of good agreement of model and reconstruction, (b) correlation between Australia
1709 and South America as an example of poor agreement. For all other combinations see Figure
1710 S9.

1711

1712

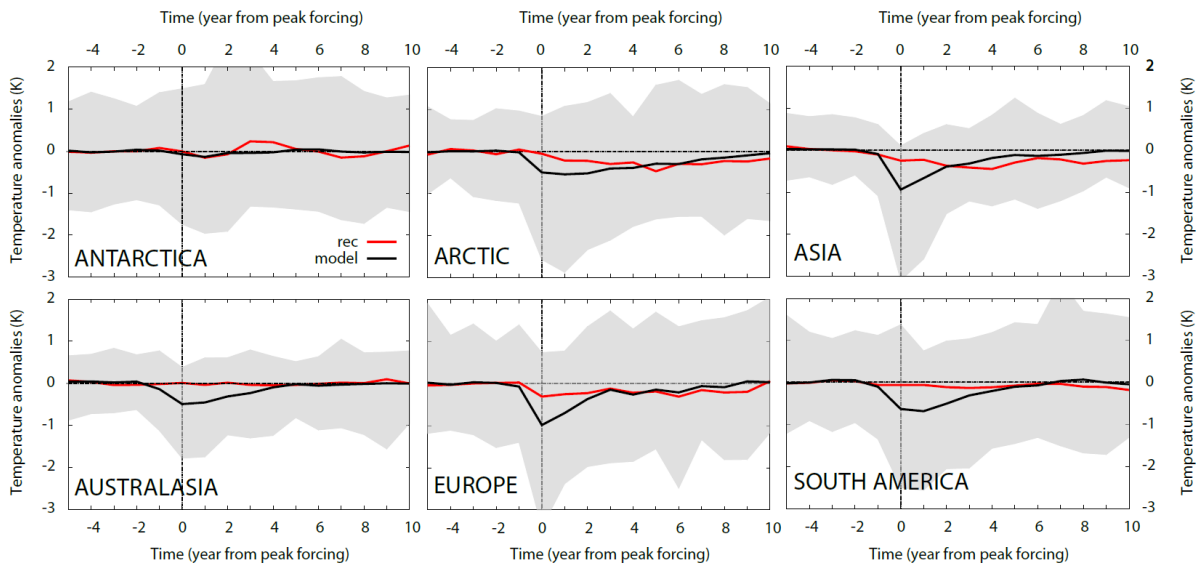
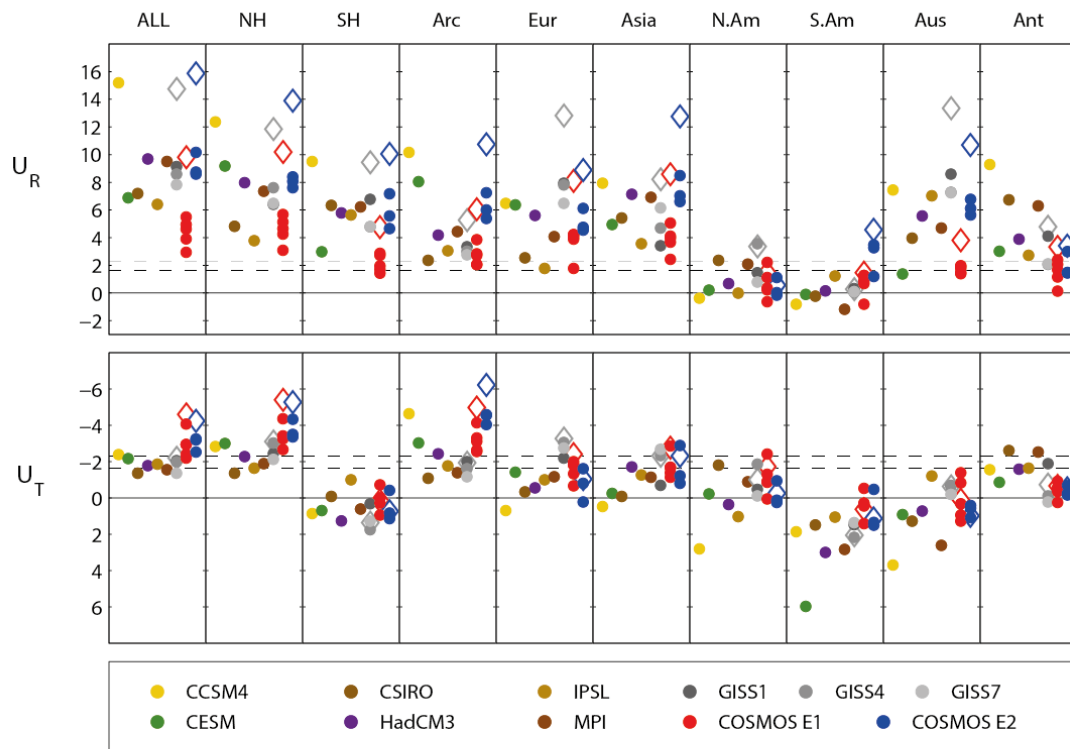


Figure 8. Superposed Epoch Analysis of the impact of the volcanic activity in the reconstructed and simulated temperatures. Superposed composites of temperature responses during selected periods when peak negative forcing in the Crowley and Unterman (2012) volcanic reconstruction are aligned. The composite is produced by selecting the 12 strongest volcanic events, starting 5 years before the date of the peak eruption and ending 10 years after the event. Each panel indicates the reconstructed (red lines) and simulated (black) composites of the temperature response for each Pages2k region. The grey shading indicates the complete range of simulated temperature responses.



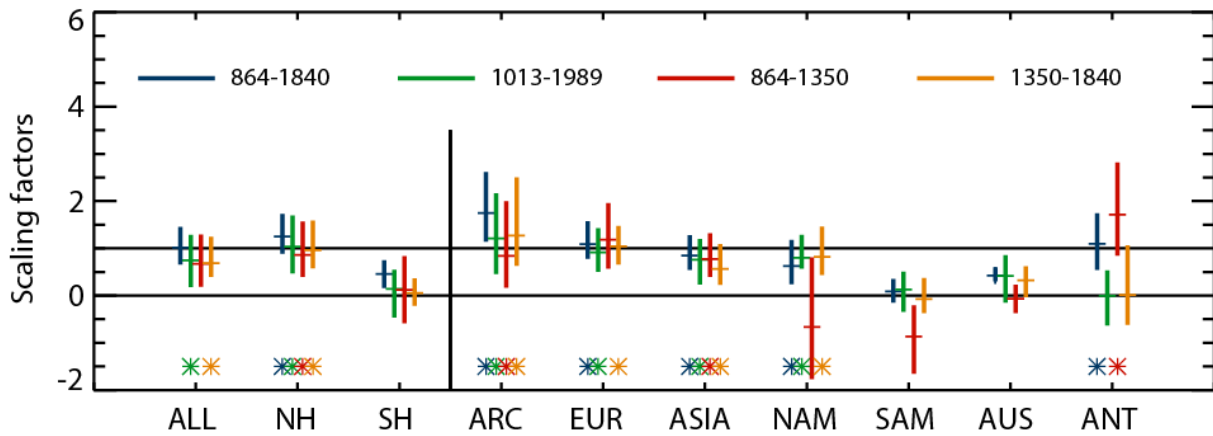
1715

1716 Figure 9: Correlation (U_R) and distance (U_T) statistics for PAGES 2k regions, with
 1717 hemispheric and global combinations of all regional data, in the period 861–1850 CE. Positive
 1718 U_R indicates that simulations and reconstructions have a positive correlation and that they
 1719 share an effect of temporal changes in external forcings. Negative U_T indicates that a forced
 1720 simulation is closer to the observed temperature variations than its own control simulation.
 1721 The analysis reveals a notably better general agreement between simulations and
 1722 reconstructions for the Northern Hemisphere as compared to the Southern Hemisphere.
 1723 Coloured dots: individual simulations. Diamonds: ensemble-mean results for COSMOS and
 1724 GISS models. Dashed lines show one-sided 5% and 1% significance levels. Note the reversed
 1725 vertical axis in the U_T graphs.

1726

1727

1728



1729

1730 Figure 10: Detection and attribution results for PAGES 2k regions. Vertical bars indicate 5-
1731 95% scaling factor ranges, with a cross indicating the best fit. Scaling factors that are
1732 significantly offset from '0' indicate that the response to forcing is detected, and those that
1733 encompass '1' indicate that the magnitude of the forced response agrees with simulations. For
1734 each region scaling ranges are shown for four different time periods (colours). For the
1735 Northern hemisphere (NH), Southern Hemisphere (SH) and global (ALL), the regressions
1736 were carried out on the combined data from all the applicable regions. An asterisk indicates
1737 that the detection analysis has been successful, namely the forced response is significantly
1738 greater than zero and that the residuals are consistent with model-based samples of internal
1739 variability.

1740

1741

1742

1743

1744

1745

1746

1747

1748

1749

1750 Table 1. Description of the model simulations

Model	# runs	Resolution	Resolution	Forcing						Reference	1751
				S	V	G	A	L	O		piControl length (yrs)
CCSM4	1	0.9° x 1.25°, L26 (atm) nominal 1°, L60 (ocn)	288 x 192, L26 (atm) 320 x 384, L60 (ocn)	10	20	30, 31, 32	40	50	60	Landrum et al. (2013)	1753 500 1754
CESM1	1	0.9° x 1.25°, L26 (atm) nominal 1°, L60 (ocn)	288 x 192, L26 (atm) 320 x 384, L60 (ocn)	11	20	30, 31, 32	40	50	1990 CE	Lehner et al. (submitted.)	1755 465 1756
CSIRO-Mk3L-1-2	1	5.63° x 3.21°, L18 (atm) 2.81° x 1.61°, L21 (ocn)	64 x 56, L18 (atm) 128 x 112, L21 (ocn)	12	21	30, 31, 32	none	none	60	Phipps et al. (2013)	1757 1150 1758 1759
GISS-E2-R	3	2° x 2.5°, L40 (atm) 1° x 1.25°, L32 (ocn)	144 x 90, L40 (atm) 288 x 180, L32 (ocn)	12	21, 20	30, 31, 32	40	50, 51	60	Schmidt et al. (2014b)	1760 1162 1761
HadCM3	1	3.75° x 2.46°, L19 (atm) 1.25° x 1.25°, L20 (ocn)	96 x 73, L19 (atm) 288 x 144, L20 (ocn)	12	21	30, 33, 32	41	51	60	Schurer et al. (2013)	1762 1199 1763
IPSL-CM5A-LR	1	3.75° x 1.88°, L17 (atm) 1.98° x 1.21°, L32 (ocn)	96 x 96, L17 (atm) 182 x 149, L32 (ocn)	10	22	30, 31, 32	none	none	60	Dufresne et al. (2013)	1764 1004 1765
MPI-ESM-P	1	1.84° x 1.84°, L47 (atm) nominal 1.5°, L40 (ocn)	196 x 98, L47 (atm) 256 x 220, L40 (ocn)	10	21	30, 31, 32	40	52	60	Jungclaus et al. (2014)	1766 1155 1767
ECHAM5/MPIOM (COSMOS)	E1:5 E2: 3	3.75° x 3.75°, L19 (atm) nominal 3°, L40 (ocn)	96 x 48, L19 (atm) 120 x 101, L40 (ocn)	13 14	21 21	32, 34 32, 34	40 40	52 52	61 61	Jungclaus et al. (2010)	1768 1000

1769 **Forcings:** S,V,G,A, L and O stands respectively for Solar, Volcanic, Greenhouse gas, Aerosols,
1770 Land use and Orbital forcing, respectively, derived from the following references:

1771 10 = Vieira and Solanki (2010) spliced to Wang et al. (2005)

1772 11 = as 10, but scaled to double the Maunder Minimum-Present Day amplitude

1773 12 = Steinhilber et al. (2009) spliced to Wang et al. (2005)

1774 13 = Krivova et al. (2007)

1775 14 = Bard et al. (2000)

1776 20 = Gao et al. (2008)

1777 21 = Crowley and Unterman (2013)

1778 22 = Ammann et al. (2007)

1779 30 = Flückiger et al. (1999, 2002); Machida et al. (1995)

1780 31 = Hansen and Sato (2004)

1781 32 = MacFarling Meure et al. (2006)

1782 33 = Johns et al. (2003)

1783 34 = CO2 diagnosed by the model.

1784 40 = Lamarque et al. (2010)

1785 41 = Johns et al. (2003)

1786 50 = Pongratz et al. (2009) spliced to Hurtt et al. (2011)

1787 51 = Kaplan et al. (2011)

1788 52 = Pongratz et al. (2008)

1789 60 = Berger (1978)

1790 61 = Bretagnon and Francou (1988)

1791

1792

1793

