

An updated assessment of near-surface temperature change from 1850: the HadCRUT5 dataset

C. P. Morice¹, J. J. Kennedy¹, N. A. Rayner¹, J. P. Winn¹, E. Hogan¹, R. E. Killick¹, R. J. H. Dunn¹, T. J. Osborn², P. D. Jones² and I. R. Simpson¹

¹ Met Office Hadley Centre, Exeter, EX1 3PB, UK

² Climatic Research Unit, School of Environmental Sciences, University of East Anglia, Norwich, NR4 7TJ, UK

Corresponding author: Colin Morice (colin.morice@metoffice.gov.uk)

Key Points:

- We have created a new version of the Met Office Hadley Centre and Climatic Research Unit global surface temperature dataset for 1850-2018.
- The new dataset better represents sparsely observed regions of the globe and incorporates an improved sea-surface temperature dataset.
- This dataset shows increased global average warming since the mid-nineteenth century and in recent years, consistent with other analyses.

Abstract

We present a new version of the Met Office Hadley Centre/Climatic Research Unit global surface temperature dataset, HadCRUT5. HadCRUT5 presents monthly average near-surface temperature anomalies, relative to the 1961-1990 period, on a regular 5° latitude by 5° longitude grid from 1850 to 2018. HadCRUT5 is a combination of sea-surface temperature measurements over the ocean from ships and buoys and near-surface air temperature measurements from weather stations over the land surface. These data have been sourced from updated compilations and the adjustments applied to mitigate the impact of changes in sea-surface temperature measurement methods have been revised. Two variants of HadCRUT5 have been produced for use in different applications. The first represents temperature anomaly data on a grid for locations where measurement data are available. The second, more spatially complete, variant uses a Gaussian process based statistical method to make better use of the available observations, extending temperature anomaly estimates into regions for which the underlying measurements are informative. Each is provided as a 200-member ensemble accompanied by additional uncertainty information. The combination of revised input datasets and statistical analysis results in greater warming of the global average over the course of the whole record. In recent years, increased warming results from an improved representation of Arctic warming and a better understanding of evolving biases in sea-surface temperature measurements from ships. These updates result in greater consistency with other independent global surface temperature datasets, despite their different approaches to dataset construction, and further increase confidence in our understanding of changes seen.

38 **Plain Language Summary**

39 We have produced a new version of a dataset that measures changes of near-surface temperature
40 across the globe from 1850 to 2018, called HadCRUT5. We have included an improved dataset
41 of sea-surface temperature, which better accounts for the effects of changes through time in how
42 measurement were made from ships and buoys at sea. We have also included an expanded
43 compilation of measurements made at weather stations on land.

44 There are two variations of HadCRUT5, produced for different uses. The first, the “HadCRUT5
45 non-infilled dataset”, maps temperature changes on a grid for locations close to where we have
46 measurements. The second, the “HadCRUT5 analysis”, extends our estimates to locations further
47 from the available measurements using a statistical technique that makes use of the spatial
48 connectedness of temperature patterns. This improves the representation of less well observed
49 regions in estimates of global, hemispheric and regional temperature change.

50 Together, these updates and improvements reveal a slightly greater rise in near-surface
51 temperature since the nineteenth century, especially in the Northern Hemisphere, which is more
52 consistent with other datasets. This increases our confidence in our understanding of global
53 surface temperature changes since the mid-nineteenth century.

54

55 **1 Introduction**

56 Observational evidence plays an essential role in our understanding of the climate, the causes of
57 the observed changes and distance travelled along predicted future trajectories. Compilations of
58 near-surface temperature measurements, as traditionally measured over land in shielded
59 enclosures and at sea by ships and buoys, as well as multi-decadal temperature records derived
60 from these compilations, are a core repository of information underpinning our understanding of
61 a changing climate. Here we present an update to one such assessment, the Met Office Hadley
62 Centre/Climatic Research Unit HadCRUT dataset (version HadCRUT.5.0.0.0, referred to
63 hereafter as HadCRUT5), incorporating additional measurements, improved understanding of
64 non-climatic effects associated with an ever-changing measurement network, and updated
65 gridding methods.

66 Global near-surface temperature analyses, based on a combination of air temperature
67 observations over land with sea-surface temperature (SST) observations, are among the longest
68 instrumental records of climate change and variability. They are routinely used in assessments of
69 the state of the climate (e.g. Blunden & Arndt, 2019). They underpin our understanding of multi-
70 decadal to centennial changes and the causes of those changes (e.g. Hartmann et al., 2013) and
71 are a key metric against which climate change policy decisions are made and progress against
72 international agreements is measured (e.g. Allen et al., 2018).

73 Analyses of multi-decadal temperature changes based on instrumental evidence are subject to
74 uncertainty. Assessments of uncertainty and the influence of non-climatic factors on observations
75 are necessary to understand the evolution of near-surface temperature throughout the
76 instrumental period. Known sources of uncertainty include spatial and temporal sampling of the
77 globe (Jones et al., 1997; Brohan et al., 2006), changes in measurement practice and
78 instrumentation (Parker 1994; Kent et al., 2017), siting of observing stations and the effects of
79 changes in their nearby environment (Parker 2006; Menne et al., 2018), and basic measurement
80 error.

81 Since the release of the predecessor of the dataset presented here, HadCRUT4 (Morice et al.,
82 2012), new analyses of near-surface temperature have been undertaken, and with them
83 understanding has improved of deficiencies in the observing network and in analysis methods.
84 This has led to updates to analyses with long pedigrees (Zhang et al., 2019; Lenssen et al., 2019),
85 the arrival of new and independent analyses (Rohde et al., 2013a; 2013b; Rohde & Hausfather,
86 2020; Yun et al., 2019), and related studies (Ilyas et al., 2017; Benestad et al., 2019; Kadow et
87 al., 2020).

88 Efforts to consolidate archives of instrumental air temperature series under the auspices of the
89 International Surface Temperature Initiative (ISTI; Rennie et al., 2014) have greatly increased
90 the availability of meteorological station series. The resulting ISTI databank underpins the
91 updated GHCNv4 air temperature data set (Menne et al., 2018) and regional subsets of station
92 series from the ISTI databank have been selectively included in updates to the CRUTEM4 and
93 CRUTEM5 datasets (Jones et al., 2012; Osborn et al., 2020). These improved data holdings have
94 increased observational coverage of regions that were previously poorly represented, including
95 the rapidly warming high northern latitudes.

96 Rohde et al. (2013a; 2013b) introduced a new land air temperature analysis developed
97 independently of pre-existing studies. This analysis included a new method for bias-adjusting
98 station records, a process that is commonly known as homogenization, and combined estimation

99 of homogenization adjustments with an independently developed spatial analysis method. The
100 study has since been extended to include analysis of HadSST3 sea-surface temperatures
101 (Kennedy et al., 2011a; 2011b) to produce a merged land-sea data product (Rohde & Hausfather,
102 2020).

103 A key uncertainty for estimating long-term change is that associated with corrections for
104 systematic errors in sea-surface temperature measurements. Comparisons of long historical SST
105 data sets (Kent et al., 2017) showed that there were differences between SST data sets which
106 were larger than the estimated uncertainties. A comparison to modern “instrumentally
107 homogeneous” data sets by Hausfather et al. (2017), found that HadSST3 (Kennedy et al., 2011a;
108 2011b) and COBE-SST-2 (Hirahara et al. 2014) underestimated recent warming. Cowtan et al.
109 (2018) compared SST products to coastal weather stations highlighting discrepancies between
110 temperature trends in land and ocean data sets. Carella et al. (2018) used characteristic daily-
111 cycles in SST measurements to infer how the measurements were made and showed that
112 previous assumptions under-estimated the prevalence of engine-room measurements.

113 Freeman et al. (2017) compiled release 3.0 of the International Comprehensive Ocean
114 Atmosphere Data Set (ICOADS) including newly digitized data. Two long-term historical SST
115 analyses, HadSST and ERSST, which are based on ICOADS, have been updated using this new
116 release. ERSST has gone through two updates – version 4 (Huang et al., 2016) and 5 (Huang et
117 al., 2017) – which extended bias adjustments to the whole SST record, implemented
118 improvements to the analysis, and quantified uncertainty. HadSST.4.0.0.0 (Kennedy et al., 2019)
119 revisited the bias adjustments applied to the data, using oceanographic measurements to
120 understand and reduce some of the key uncertainties in HadSST3.

121 Recent updates to instrumental near-surface temperature data products have brought
122 improvements in their assessment of uncertainty, and in provision of uncertainty information for
123 use in onward analyses. Ensemble uncertainty assessments have become commonplace in air
124 temperature datasets (Morice et al., 2012; Menne et al., 2018) and sea-surface temperature
125 datasets (Kennedy et al., 2011b; Huang et al., 2016; Huang et al., 2019; Kennedy et al., 2019).

126 The NOAA GlobalTemp version 5 analysis (Zhang et al., 2019; Huang et al., 2019) updates
127 previous NOAA analyses (Smith et al., 2008) by bringing together updates to underpinning data
128 holdings over land (Menne et al., 2018) and merges the expanded land data holdings of GHCNv4
129 with the updated ERSSTv5 data set. An ensemble uncertainty assessment is included (Huang et
130 al., 2019), sampling the uncertainty in parametric choices in the SST adjustments procedure, the
131 station series homogenization algorithm (Menne et al., 2018) and the spatial analysis method
132 used.

133 The NASA Goddard Institute for Space Studies GISTEMPv4 analysis (Lenssen et al., 2019)
134 introduces an updated uncertainty assessment, applying the GISTEMP spatial analysis methods
135 to the 100-member GHCNv4 ensemble of homogenized station series and basing SST
136 uncertainty assessments on the ERSSTv4 ensemble. Additional uncertainty associated with the
137 production of spatial analyses from incomplete station data is assessed by sub-sampling
138 reanalysis fields from a selection of modern reanalyses.

139 Coverage of instrumental records of near-surface temperature changes is characterized by often
140 sparse and non-uniform sampling of the globe. Assessments of uncertainty in global and regional
141 average temperature changes have found that sparse data coverage is the most prominent source
142 of uncertainty over monthly to decadal timescales (Brohan et al., 2006; Morice et al., 2012),

143 outweighing uncertainty arising from changes in observing methods. Recent studies have also
144 shown that poor representation of some regions, notably the rapidly warming high northern
145 latitudes, may have contributed to an underestimation of globally averaged temperature changes
146 in recent years (Cowtan and Way, 2014; Karl et al., 2015).

147 While efforts have been made to increase data coverage in the CRUTEM4 and now CRUTEM5
148 data set through inclusion of additional meteorological station data in less well-observed regions
149 (Jones et al., 2012; Osborn et al., 2020) and marine data holdings expanded to include recently
150 digitized marine reports (Freeman et al. 2017), statistical analysis methods were not used in
151 HadCRUT4 or its underpinning land and marine datasets to infer temperature changes in regions
152 where measurements are not available. An independent application of local statistical
153 interpolation methods to HadCRUT4, in a study by Cowtan and Way (2014), found that
154 statistically infilled reconstructions showed recent warming over high latitude regions that is not
155 proportionately represented in global mean temperatures calculated from the non-infilled
156 HadCRUT4 data set. The study also included an analysis that used satellite-based upper air
157 temperature estimates as a proxy for near-surface temperature variability in the gaps in data
158 coverage in HadCRUT4, which also showed warming in these high latitude regions. This high-
159 latitude signal contributed to an increase in the assessed rate of change of global average
160 temperatures since the beginning of the 21st century.

161 Unlike HadCRUT4, other existing near-surface temperature datasets utilize statistical analysis
162 methods to infer spatial fields from scattered observations. Analysis methods based on spatial
163 covariance structure, known variously as optimal interpolation (e.g. as used in Reynolds &
164 Smith, 1994), kriging (e.g. as used in Cowtan & Way, 2014), Gaussian process regression
165 (Rasmussen & Williams, 2006) and variants thereof, have a long history of use, particularly in
166 analyses of sea-surface temperatures (Reynolds et al., 2002; Reynolds & Smith, 1994; Donlon et
167 al. 2012). These methods use knowledge of the covariance structure of spatial fields to infer field
168 values as weighted averages of observations in locations with strong covariation. Typically,
169 weighting is based on a statistical model in which nearby locations are expected to covary
170 strongly and distant locations weakly. Methods of this form are a core part of the Rohde &
171 Hausfather (2020) analysis and of the analysis of Cowtan and Way (2014). The GISTEMP data
172 set also uses a distance-weighted average that, while similarly applying a weighted average of
173 local observations, does not make use of a covariance model and so does not classify as a kriging
174 type analysis.

175 A second form of spatial analysis methods that are commonly applied in instrumental climate
176 analyses, reduced space methods, decompose spatial temperature variability into a finite,
177 typically orthogonal, set of spatial patterns of variability (Kaplan et al., 1997). These patterns are
178 generally, but not necessarily, global in extent. Spatial reconstructions are then formed as a
179 weighted sum of these patterns. The Empirical Orthogonal Teleconnection (Smith et al., 2008)
180 method employed within the NOAA GlobalTemp v5 analysis falls within this category of
181 reduced space algorithms, employing a finite set of locally defined spatial patterns that are fit to
182 the available data.

183 A recent assessment of the use of neural networks to estimate missing values in the HadCRUT4
184 dataset (Kadow et al., 2020) expands the ensemble of methods used to reconstruct global
185 temperatures. Derived global temperature series show good agreement with prior studies using
186 more traditional methods.

187 Traditionally, surface temperature data sets have combined air temperatures over land with sea-
188 surface temperatures over the ocean, rather than the more natural choice of air temperatures over
189 the ocean. SST measurements are currently far more numerous than marine air temperature
190 (MAT) measurements because SST can be readily measured by automatic sensors on drifting
191 buoys as well as being retrieved from satellite measurements of radiances, while observational
192 sampling of MAT has been in recent decline (Berry & Kent., 2017). There are significant biases
193 in daytime marine air temperature observations (Berry et al., 2004). Night-time measurements
194 have therefore been used to develop observational records of marine air temperature changes
195 (Kent et al. 2013), with up-to-date independent assessments of historical night-time MAT
196 becoming available only recently (Junod & Christy 2020, Cornes et al., 2020). Anomalies in
197 MAT and SST have been expected to be similar over long space and time scales due to the
198 strong physical link between the two. However, Cowtan et al. (2015) showed that MAT and SST
199 changes simulated in coupled climate models differ, with MAT warming slightly faster than
200 SST, affecting comparisons of observed and modelled global temperature change if care is not
201 taken to ensure an “apples to apples” comparison. They also found that decisions about how to
202 handle marginal sea-ice areas could affect the estimated changes, depending on the use of SST or
203 MAT. Therefore, while there is good motivation for the use of MAT (Cowtan et al., 2015;
204 Richardson et al., 2016), there are currently challenges relating to the MAT observational
205 network (Berry & Kent, 2017) that provide an observational rationale for the continued use of
206 SST in monitoring global surface temperature variability and change until these challenges are
207 addressed.

208 Recent developments in satellite retrievals of surface skin temperatures present a new possibility
209 for near-surface temperature monitoring, bringing the potential for detailed spatial information
210 with sustained measurement over a time frame that is now of sufficient length for climate
211 studies. Recent work (Rayner et al., 2020) has explored the potential of combining air
212 temperature information inferred from satellite skin temperatures with traditional *in situ*
213 observations, expanding on the understanding of relationships between satellite-derived skin
214 temperatures and traditional near-surface air temperature observations, and on the stability of
215 these relationships over time that is required to construct merged data products. Alternatively,
216 dynamical reanalyses, that combine numerical weather prediction models with a range of varied
217 observational data sources, are increasingly being used to monitor the climate (e.g. ERA5,
218 Hersbach et al., 2020; JRA-55, Kobayashi et al., 2015; and MERRA-2, Gelaro et al., 2017).
219 These alternative sources of near-surface temperature data provide useful information in
220 locations that are not well represented in traditional near-surface temperature datasets. However,
221 in all cases, understanding of non-climatic effects affecting observations and arising from
222 analysis methods is required when combining observations from multiple sources.

223 Here, two ensemble surface temperature datasets are presented. The first, the “HadCRUT5 non-
224 infilled dataset”, adopts the gridding and ensemble generation methods of HadCRUT4 (Morice
225 et al., 2012). The second, the “HadCRUT5 analysis”, uses a statistical infilling method to
226 improve the representation of sparsely observed regions. Through application of the statistical
227 infilling method to the HadCRUT5 non-infilled ensemble, the HadCRUT5 analysis ensemble
228 samples the uncertainty in the gridded near-surface temperature data that arises from residual
229 biases in observational data after correction, for example associated with uncertainty in changes
230 in instrumentation and measurement practices at meteorological stations (Brohan et al., 2006;
231 Morice et al., 2012) and changes in sea-surface temperature measurement methods (Kennedy et
232 al., 2019). It also samples the effects of basic measurement uncertainty, uncertainty arising from

233 estimation of gridded temperature fields from a finite number of observations and residual
234 uncertainties associated with individual marine measurement platforms, where information
235 identifying individual platforms is available (Kennedy et al., 2019). Statistical reconstruction
236 uncertainty is also encoded in the HadCRUT5 analysis ensemble, producing an ensemble that
237 samples a greater range of sources of uncertainty than was possible in HadCRUT4. Thus, the
238 new ensemble analysis communicates the major known sources of uncertainty in an easily
239 accessible way.

240 The remaining sections of this paper are structured as follows. Section 2 describes the data sets
241 used as inputs and for comparison. Section 3 provides an overview of the methods used to
242 construct HadCRUT5. Results are presented in Section 4 with conclusions and discussion in
243 Section 5.

244 **2 Input Datasets**

245 2.1 HadSST.4.0.0.0

246 Version 4 of the Met Office Hadley Centre Sea-Surface Temperature data set, HadSST.4.0.0.0
247 (Kennedy et al., 2019), is based on *in situ* measurements of SST from ships and buoys. The ship
248 and buoy measurements are taken from ICOADS release 3.0 (Freeman et al. 2017) from 1850 to
249 2014 and release 3.0.1 from 2015 to 2018. From 2016 onwards, measurements from drifting
250 buoys are taken from the Copernicus Marine Environment Monitoring Service, as buoy data in
251 ICOADS were incomplete following a change in data-transmission codes in late 2016. Early
252 measurements made using buckets are adjusted using a physically based model of heat lost from
253 water-sampling buckets (Folland and Parker 1995; Rayner et al., 2006). From the 1940s
254 onwards, ship measurements are adjusted based first on comparisons with near-surface
255 oceanographic measurements (Atkinson et al., 2014) and then, from the early 1990s onwards, on
256 comparisons with buoy measurements. The resulting HadSST.4.0.0.0 data set is presented as
257 anomalies relative to 1961-1990 on a 5° latitude by 5° longitude grid and is representative of
258 SST as measured by drifting buoys at an approximate depth of 20 cm.

259
260 Overall, the global SST change estimated from HadSST.4.0.0.0 is larger than that estimated from
261 HadSST.3.1.1.0 (and earlier versions). This is due to two factors. First, new estimates of biases
262 associated with measurements made in ships' engine rooms show that these biases have declined
263 since the 1950/60s, artificially reducing the long-term change represented in the underlying data
264 and in earlier versions of HadSST. Second, a greater proportion of measurements during the
265 1961-1990 period were estimated to have been made in ships' engine rooms. Other changes
266 include: using buoys as a reference data set; producing ensemble members with step changes in
267 the time evolution of the proportions of canvas and wooden buckets in the early 20th century
268 alongside ensemble members which assume a linear transition; estimating the fraction of
269 incorrect metadata using comparisons with oceanographic measurements; and using comparisons
270 with oceanographic measurements to narrow the range of plausible transition dates from canvas
271 buckets to modern rubber buckets (see Kennedy et al. (2019) for a detailed discussion).

272
273 Uncertainty in HadSST.4.0.0.0 is split into three main components associated with: pervasive
274 systematic errors; systematic errors from individual ships or buoys; and uncorrelated errors from
275 individual measurements and incomplete grid-box sampling. The pervasive systematic errors,
276 which have complex temporal and spatial correlations, are represented using a 200-member

277 ensemble generated by varying uncertain parameters and choices in the bias adjustment scheme.
278 The systematic errors are represented using covariance matrices that encode the error
279 covariances between grid cells that arise from ships making measurements in multiple grid cells
280 in a month. Finally, uncertainties from uncorrelated errors are provided as gridded fields. Note
281 that these uncertainty components do not span the full range of uncertainty. In particular,
282 structural uncertainty remains (Thorne et al., 2011) and there may be an underestimate in the
283 systematic error component because it does not currently deal explicitly with errors that correlate
284 at the level of national shipping fleets (Chan & Huybers, 2019) or with marine reports that lack
285 ship call signs or other identifying information (Carella et al., 2017).

286

287 2.2 CRUTEM.5.0.0.0

288 Monthly averages of near-surface air temperature measured at weather stations over the land
289 surface for 1850-2018 are obtained from CRUTEM.5.0.0.0 (Osborn et al., 2020, referred to
290 hereafter as CRUTEM5). The CRUTEM station database is a collection of station series obtained
291 from National Meteorological and Hydrological Services (NMHSs) and large collections such as
292 the European Climate Assessment and Dataset (Klein Tank et al., 2002). CRUTEM incorporates
293 corrections that NMHSs apply to their own data to minimize the impact of changes in weather
294 station instrumentation or location on the measurement series. The monthly average temperatures
295 from stations are subjected to quality control, converted to anomalies (differences from their
296 1961–1990 means) and then averaged into 5° latitude by 5° longitude grid boxes.

297

298 CRUTEM5 has improved quality control checks that: (i) improve the flagging of incorrect data
299 during 1941-1990; (ii) reduce the trend towards increased flagging of suspect data outside of the
300 1941-1990 period; and (iii) reduce the number of genuine extreme values that are erroneously
301 flagged as incorrect, e.g. during coherent extreme events such as summer 2003 in Europe (see
302 Osborn et al. (2020) for details). The station database has been expanded such that the number of
303 those stations with sufficient data to estimate temperature anomalies has grown from 4842 in
304 CRUTEM.4.0.0.0 (as used in Morice et al., 2012) to 7983 in CRUTEM5 (Osborn et al., 2020).
305 Most of the new data acquisitions are in already-sampled regions, so the number of grid-box
306 values is only moderately expanded (by 9%) relative to CRUTEM.4.0.0.0.

307

308 The changes in temperature seen in hemispheric or global averages since 1850 are not sensitive
309 to these updates, but some regional differences are apparent. Osborn et al. (2020) describes the
310 effects of updates since CRUTEM.4.0.0.0, and of updates since the more recent
311 CRUTEM.4.6.0.0 (as used in HadCRUT.4.6.0.0), in detail.

312 An alternative gridding method was explored in Osborn et al. (2020) for CRUTEM5 to address
313 the under-representation of high latitude stations in the standard gridding. This alternative
314 method allows each station to contribute to more than one neighboring 5° latitude by 5°
315 longitude grid box on the same latitude, where the number of grid boxes to which each station
316 can contribute is determined by an inverse cosine latitude weighting. In the current paper, the
317 alternative gridding method is not used because (a) the uncertainty model for the CRUTEM5
318 grids, as documented in Brohan et al. (2006), only applies to the standard gridding approach
319 (where each station contributes to only one grid box); and (b) the issue of high-latitude sampling
320 is dealt with here by statistical infilling.

321 HadCRUT5 uses an ensemble version of the CRUTEM5 uncertainty model. The HadCRUT5
322 non-infilled ensemble grids and accompanying uncertainty grids are produced from the
323 CRUTEM5 station temperature anomaly series, following the methods of Morice et al. (2012), as
324 described in Section 3.2.

325 2.3 HadISST.2.2.0.0

326 We use sea ice concentration from the Met Office Hadley Centre sea-Ice and Sea Surface
327 Temperature data set, HadISST.2.2.0.0 (an update to Titchner and Rayner (2014)), on a 1°
328 latitude by 1° longitude grid to determine the presence or absence of sea ice in any individual
329 ocean grid box in each month from 1850 to 2018.

330 HadISST.2.2.0.0 is updated relative to version 2.1.0.0 in the following ways: (i) reinstatement of
331 a small number of erroneously-removed sea-ice-filled grid boxes after 1978; (ii) an alteration to
332 the adjustments applied to the National Ice Center charts (used to determine the ice edge between
333 1972 and 1978) correcting a low-bias in the HadISST.2.1.0.0 fields in the Arctic then; and (iii)
334 an improvement in the interpolation applied between two atlas-derived climatologies used to
335 determine ice extents in the Antarctic to produce a smoother transition between them and
336 between 1962 and the start of monthly observations in 1972.

337 2.4 ERA5

338 We have used monthly ERA5 analysis 2 m air temperature data from 1979-2018 (Hersbach et
339 al., 2020) for coverage uncertainty estimation and for comparison of global and regional
340 diagnostics. ERA5 was produced using 4D-Var data assimilation in the European Centre for
341 Medium-range Weather Forecasts' (ECMWF) Integrated Forecast System (IFS). We used the
342 (31 km) high resolution realization.

343 2.5 Other comparison data

344 Four comparison data sets are used here: NOAAGlobalTemp version 5 (Zhang et al., 2019;
345 Huang et al., 2019), GISTEMP version 4 (Hansen et al., 2010; Lenssen et al., 2019), Berkeley
346 Earth (Rohde & Hausfather, 2020) and Cowtan & Way (Cowtan & Way, 2014).

347 NOAAGlobalTemp version 5 is based on the Global Historical Climatology Network (GHCN)
348 version 4 land station data set (Menne et al., 2018) and the Extended Reconstruction Sea Surface
349 Temperature (ERSST) data set version 5 (Huang et al., 2017). Station records in GHCN v4 are
350 homogenized using an automated algorithm. SSTs are adjusted using comparisons with marine

351 air temperature and latterly drifting buoys. The data are interpolated using Empirical Orthogonal
352 Teleconnections, providing improved coverage, although coverage does not extend fully into the
353 polar regions.

354 GISTEMP version 4, like NOAAGlobalTemp v5, is based on a combination of GHCN v4 and
355 ERSST v5. The SST data are interpolated as in NOAAGlobalTemp. Land surface air
356 temperatures are interpolated from station data within a 1200km radius. Extrapolated land
357 surface air temperatures are used over the oceans in sea-ice covered areas. Coverage of the
358 GISTEMP data set is quasi-global in the past twenty years, with good coverage of the poles and
359 other data-sparse regions from interpolated station data.

360 The Berkeley Earth data set (Rohde & Hausfather, 2020) uses a kriging-based technique to
361 interpolate and homogenize station data. A kriging based technique is also applied to SSTs from
362 the HadSST3 data set to provide coverage over the whole globe. The version of the data set that
363 uses extrapolated land-surface air temperatures over the oceans in sea-ice covered areas is used
364 here.

365 Cowtan and Way (2014) is based on the HadCRUT4 data set. The land and ocean data are
366 interpolated using kriging. Grid cells that contain data in HadCRUT4 are not modified during
367 interpolation (in contrast to the kriging of HadSST3 data in the Berkeley Earth data set). As with
368 GISTEMP and Berkeley Earth, extrapolated land-surface air temperatures are used over the
369 oceans in sea-ice covered areas.

370 The Berkeley Earth (1° latitude by 1° longitude resolution) and ERA5 (0.25° latitude by 0.25°
371 longitude resolution) analyses were regridded to 5° latitude by 5° resolution using an area-
372 weighted average of all grid cells falling within a HadCRUT5 5° grid cell. Cowtan and Way and
373 NOAAGlobalTemp were obtained on a 5° grid. The GISTEMP data, which were obtained on a
374 2° grid, were not regridded.

375

376 **3 Methods**

377 Two gridded datasets are provided as part of HadCRUT5. The first version of the dataset is
378 produced without statistical infilling, referred to here as the “HadCRUT5 non-infilled dataset”,
379 following the methods of Morice et al. (2012), and is intended for use in applications where
380 statistical infilling is not desired. This is accompanied by a second version of the dataset,
381 hereafter referred to as the “HadCRUT5 analysis”, that is produced using a statistical method to
382 estimate more-complete temperature anomaly fields.

383 The HadCRUT5 non-infilled dataset and the HadCRUT5 analysis are produced in the following
384 steps. First, an ensemble land-surface air temperature dataset, with accompanying additional
385 uncertainty information, is generated from the CRUTEM5 station data (Section 3.2). The land
386 dataset is then merged with sea-surface temperature anomaly information from HadSST4
387 through a weighting method based on the land area fraction (Section 3.4) to produce the non-
388 infilled dataset. Next, monthly fields are estimated separately for the land surface air temperature
389 dataset and for HadSST4 using a statistical method to create an ensemble analysis for each
390 (Section 3.3). The separate land and ocean analyses are then merged into a combined land and

391 ocean ensemble analysis using a land-sea weighting scheme that also accounts for sea ice
 392 coverage (Section 3.4). Global and regional time series are then computed from the two merged
 393 datasets, following the methods of Morice et al. (2012) with updates to the method used to
 394 estimate uncertainty associated with incomplete observational coverage described in Section 3.5.
 395 Error models for each dataset are described in Section 3.1. Full details of uncertainty propagation
 396 for land and ocean merging and global and regional time series are provided in the Supporting
 397 Information.

398

399 3.1 The HadCRUT5 error models

400 This section outlines the terms of the error model for grids and time series of the HadCRUT5
 401 non-infilled dataset and the HadCRUT5 analysis. Further details are given in the Supporting
 402 Information.

403 The error models are split into components according to the way that uncertainty information is
 404 presented in HadCRUT5. The sources of uncertainty modelled in HadCRUT5 are grouped
 405 according to their correlation structure to allow uncertainties to be propagated appropriately into
 406 derived diagnostics such as regional average time series.

407 3.1.1 The HadCRUT5 non-infilled dataset

408 The error model for the non-infilled dataset describes the estimate of temperature anomaly $\hat{T}(s, t)$
 409 at spatial location s and time t as a sum of the true temperature anomaly $T(s, t)$ and three error
 410 terms: a bias term $\varepsilon_b(s, t)$ representing biases with large-scale spatial and temporal structure; a
 411 partially correlated error term $\varepsilon_p(s, t)$ for errors with typically local structure; and an
 412 uncorrelated error term $\varepsilon_u(s, t)$ describing errors that are independent between spatial and
 413 temporal locations. The full error model for non-infilled fields is given by:

414

$$\hat{T}(s, t) = T(s, t) + \varepsilon_b(s, t) + \varepsilon_p(s, t) + \varepsilon_u(s, t) \quad (1)$$

415

416 This error model for the merged dataset matches the structure of the error model for the land
 417 dataset and for HadSST4. For the land dataset, the contributions to the bias term are the land
 418 station homogenization error, urbanization and biases from non-standard measurement
 419 enclosures. There is no contribution to the partially correlated term and the uncorrelated term
 420 models the within grid box measurement and sampling uncertainties (Morice et al., 2012). For
 421 HadSST4, the bias term models the effects of residual errors in the adjustments applied to
 422 account for changes in measurement methods, the partially correlated term models the effects of
 423 residual biases associated with individual observing platforms, and the uncorrelated term models
 424 the within grid cell measurement and sampling uncertainties (Kennedy et al., 2019).

425 The HadCRUT5 non-infilled ensemble samples the uncertainties for the combination $T(s, t) +$
 426 $\varepsilon_b(s, t)$. The uncertainties for partially correlated and uncorrelated errors are not encoded into

427 the non-infilled ensemble. Instead, uncertainty information for partially correlated errors $\varepsilon_p(s, t)$
 428 is provided in spatial error covariance matrices and uncertainties for uncorrelated errors $\varepsilon_u(s, t)$
 429 are provided for each observed grid cell.

430 The error model for estimates of spatial average time series $\hat{T}(t)$ derived from the gridded data is
 431 given as a sum of the true temperature anomaly time series $T(t)$ and four error terms:

432

$$\hat{T}(t) = T(t) + \varepsilon_b(t) + \varepsilon_p(t) + \varepsilon_u(t) + \varepsilon_c(t) \quad (2)$$

433

434 Here $\varepsilon_b(t)$ is the effect of the bias term propagated into the spatial average, $\varepsilon_p(t)$ is the effect of
 435 the partially correlated term, $\varepsilon_u(t)$ the effect of the uncorrelated error term. The fourth error
 436 term, $\varepsilon_c(t)$, is the error in estimating the spatial average from incomplete spatial coverage, with
 437 missing grid cells resulting from limitations in the spatial sampling provided by the observation
 438 network. Full details of uncertainty propagation for each of these terms are given in the
 439 Supporting Information.

440 3.1.2 The HadCRUT5 analysis

441 An overview of the HadCRUT5 analysis is provided in Section 3.4 and a detailed description of
 442 methods is provided in Appendix A. The HadCRUT5 analysis error model has fewer terms than
 443 that of the non-infilled dataset as the analysis methods combine multiple sources of error into a
 444 single analysis error term. The error model for the HadCRUT5 analysis defines the temperature
 445 anomaly estimate as the sum of the true temperature $T(s, t)$ and the analysis error $\varepsilon_a(s, t)$:

446

$$\hat{T}(s, t) = T(s, t) + \varepsilon_a(s, t) \quad (3)$$

447

448 The analysis error term combines all errors that are modelled in the Gaussian process analysis,
 449 both spatial reconstruction errors and observational errors, as described in Appendix A. The
 450 analysis ensemble samples the analysis uncertainty such that each ensemble member is a sample
 451 of $T(s, t) + \varepsilon_a(s, t)$.

452

453 For the HadCRUT5 analysis, errors in global and regional average time series are derived as a
 454 combination of the propagated analysis error and $\varepsilon_a(t)$ and an additional coverage error term
 455 $\varepsilon_c(t)$ that represents the error in estimating the spatial average from incomplete analysis grids,
 456 noting that this coverage error term differs from that of the non-infilled dataset due to the
 457 different spatial coverage of the analysis.

458

$$\hat{T}(t) = T(t) + \varepsilon_a(t) + \varepsilon_c(t) \quad (4)$$

459

460 The propagation of uncertainty associated with these errors is described in the Supporting
461 Information.

462 3.2 Ensemble land air temperature data set

463 As in the previous versions of HadCRUT, near-surface air temperature information over land is
464 derived from the CRUTEM data set. As in Morice et al. (2012), an ensemble air temperature data
465 set is produced by sampling from the distributions of known uncertainty in station temperature
466 records. The station data on which the ensemble grids are based has been updated to now use the
467 CRUTEM.5.0.0.0 data set (Osborn et al., 2020).

468 A detailed description of the land air temperature ensemble sampling method can be found in
469 Morice et al. (2012). The sampling approach is designed so that the effects of known sources of
470 residual systematic error in station anomaly series can be quantified for regional statistics and
471 time series. The ensemble size has been increased to 200 members for HadCRUT5 to match the
472 200-member HadSST4 ensemble.

473 The sampling method is as follows. Samples are drawn from the distributions of known
474 uncertainties during the station gridding process. Residual homogenization error and uncertainty
475 in climatology normal information are sampled from distributions described in Brohan et al.
476 (2006) and encoded into realizations of individual station series prior to gridding. The systematic
477 effects of residual urbanization errors (Brohan et al., 2006; Parker, 2010) and non-standard
478 sensor enclosures (Parker, 1994; Folland et al., 2001) are sampled and encoded into the gridded
479 ensemble at a regional level, again following the method of Morice et al. (2012).

480 Additional uncertainty information for errors that are uncorrelated between grid cells (e.g. from
481 measurement error or incomplete sampling of a grid cell) is not encoded into the land ensemble.
482 Instead, these measurement and sampling-related uncertainties are provided as additional
483 uncertainty information outside of the ensemble, as in Morice et al. (2012).

484

485 3.3 Spatial analysis of temperature anomaly fields

486 HadCRUT5 now includes an ensemble spatial analysis that reconstructs more spatially extensive
487 anomaly fields from the available observational coverage. The purpose of this analysis is to: (1)
488 reduce uncertainty and bias associated with estimation of global and regional climate diagnostics
489 from incomplete and uneven observational sampling of the globe; (2) provide improved
490 estimates of temperature fields in all regions; and (3) provide a method to quantify uncertainty in
491 anomaly patterns.

492 We adopt a Gaussian process based method for spatial analysis that is closely related to the
493 ordinary kriging approach (Rasmussen & Williams, 2006), and apply the method independently
494 to land air temperature and sea-surface temperature observations before merging the two to
495 produce a global analysis. The method models monthly temperature anomaly fields as
496 realizations of a Gaussian processes with a simple covariance structure, defined as a function of
497 the distance between locations, and an *a priori* unknown mean, and accounts for observational
498 uncertainty. A detailed description of the analysis method is presented in Appendix A.

499 The Gaussian process method is applied to the 5° latitude by 5° longitude gridded anomaly fields
 500 of the land ensemble and the HadSST4 ensemble. The additional observational uncertainty terms
 501 that accompany these input ensembles are provided to the Gaussian process estimation as
 502 monthly error covariance matrices. The spatial reconstructions are based upon a model of the
 503 covariance structure of the 5° latitude by 5° longitude anomaly fields. This covariance structure
 504 is modelled using a Matérn covariance function, for which the covariance decays as a function of
 505 Euclidian distance between locations. The parameters of the Matérn covariance function are
 506 fitted separately for land air temperature and sea-surface temperature anomalies (see Appendix
 507 A.2), representing typical variability in each domain.

508 As a Bayesian method, the approach provides a framework for assessing analysis uncertainties
 509 and provides the capability to draw samples from the posterior distribution of the analysis. We
 510 generate an ensemble of field estimates through application of the analysis method to each input
 511 ensemble member and then drawing samples from the posterior distributions of the Gaussian
 512 process estimates. The land and ocean analysis ensembles combine all sources of uncertainty
 513 represented in the input gridded datasets whilst respecting the estimated covariance structure of
 514 the temperature anomaly field so that each ensemble member is a plausible spatial analysis of the
 515 temperature anomaly field.

516 The analysis has limited capability to reconstruct temperatures at long distances from available
 517 observations, as the field estimates are based on a model of local covariance structure. We
 518 therefore introduce criteria for excluding regions where there is not a strong observational
 519 constraint on the analysis (see Appendix A.4). The masked land air temperature and sea-surface
 520 temperature anomaly ensembles are then merged, as described in Section 3.4.

521 3.4 Blending land air temperatures with sea-surface temperature data

522 The 200-member ensemble land air temperature data set based on CRUTEM5 and the 200-
 523 member HadSST4 are merged as a weighted average of the 5° latitude by 5° longitude land and
 524 marine fields. Two versions of the data set are provided: one that uses the spatial analysis
 525 method presented in Section 3.2 and one that does not.

526 3.4.1 Merging non-infilled datasets

527 For the non-infilled dataset, the land air temperature ensemble and HadSST4 ensemble members
 528 are merged following the methods of Morice et al. (2012). The temperature anomaly $T(s, t)$ at
 529 location s and time t is defined as the weighted average of the air temperature anomaly $T_L(s, t)$
 530 and sea surface temperature anomaly $T_M(s, t)$, with weights $f(s, t)$:

$$531 \quad T(s, t) = f(s, t)T_L(s, t) + (1 - f(s, t))T_M(s, t) \quad (5)$$

532

533 The weighted average is based on the areal fraction of land and sea in a 5° latitude by 5°
 534 longitude grid cell using the same land fraction data set as HadCRUT4, originally derived from
 535 the Operational Sea Surface Temperature and Sea Ice Analysis (OSTIA; Donlon et al., 2012)
 536 0.05° land mask information. As in HadCRUT4, land air temperature information receives a

537 minimum weighting of 25% to prevent island stations from receiving near-zero weighting.
538 Where only one of the land air temperature or sea-surface temperature data sets are available, the
539 available data source receives 100% weighting.

540 Methods for merging the uncertainty fields and measurement error covariance information for
541 land and marine data sets are unchanged from those described in Morice et al. (2012) and are
542 detailed in the Supporting Information.

543 3.4.2 Merging land and ocean analyses

544 The land-sea weighting scheme is modified for the HadCRUT5 analysis. Areas of sea ice are
545 treated as if they were land in the weighting (consistent with the approach used by Cowtan &
546 Way (2014)), so that temperature anomalies over sea ice are reconstructed as part of the air
547 temperature analysis rather than the SST analysis.

548 Sea ice concentrations are obtained from the HadISST.2.2.0.0 data set. Where the ice
549 concentration on the native 1° latitude by 1° longitude HadISST.2.2.0.0 grid exceeds 15%, the
550 threshold value used to define the ice edge in Titchner and Rayner (2014), the area is considered
551 to be ice covered for the purpose of deriving weights. Ice concentrations below 15% are treated
552 as open water. For each HadISST.2.2.0.0 grid cell, a value of one is set if the sea-ice
553 concentration is greater than 15% and zero otherwise. On the 5° latitude by 5° longitude
554 HadCRUT5 grid, the fractional area of water covered by sea ice is then obtained through area-
555 weighted averages of the non-land 1° grid cells of ones and zeroes. This area of ice-covered
556 water is treated as land when deriving weights for land and ocean analyses.

557 The 25% minimum weighting for land air temperature is retained for any 5° latitude by 5°
558 longitude grid cells that are observed in the non-infilled land air temperature data set so that
559 information from island stations is not lost in the averaging. This minimum weighting is not
560 applied in grid cells that are not directly observed. Reconstructed land air temperatures are not
561 used over 100% sea regions where there are no land stations or sea ice and, similarly,
562 interpolated SST is not used over 100% land regions. This prevents extrapolation of land air
563 temperature far into ocean regions and prevents inland extrapolation of SSTs.

564 3.5 Estimating uncertainty arising from incomplete coverage

565 Spatial fields of temperature anomalies in the non-infilled HadCRUT5 data set and the
566 HadCRUT5 analysis are not globally complete. Variability in regions of the world that are not
567 represented in the spatial fields gives rise to uncertainty in global and regional time series. For
568 the non-infilled HadCRUT5, the coverage uncertainty accounts for regions of the globe where
569 insufficient observations are available to compute grid cell average anomalies in the underlying
570 air temperature and SST data sets. For the HadCRUT5 analysis, the coverage uncertainty
571 accounts for the masked regions of the analysis that are not well constrained by observations.

572 Coverage uncertainty is assessed by sub-sampling globally-complete reanalysis fields to the
573 coverage of HadCRUT5 using the method presented in Brohan et al. (2006) and Morice et al.
574 (2012), which is described in detail in the Supporting Information. The approach is updated here
575 to use the recently-released ERA5 reanalysis (Hersbach et al., 2020) as the globally-complete
576 reference data set, in place of the previously used NCEP/NCAR reanalysis (Kalnay et al., 1996).

577 Temperature anomalies are computed for the ERA5 monthly 2 m air temperature grids,
578 referenced to the period of ERA5 that overlaps with our climatology period: 1979-1990.
579 Anomalies are then averaged to the 5° latitude by 5° longitude grid used in HadCRUT5 to
580 produce the reference fields for the coverage uncertainty calculations.

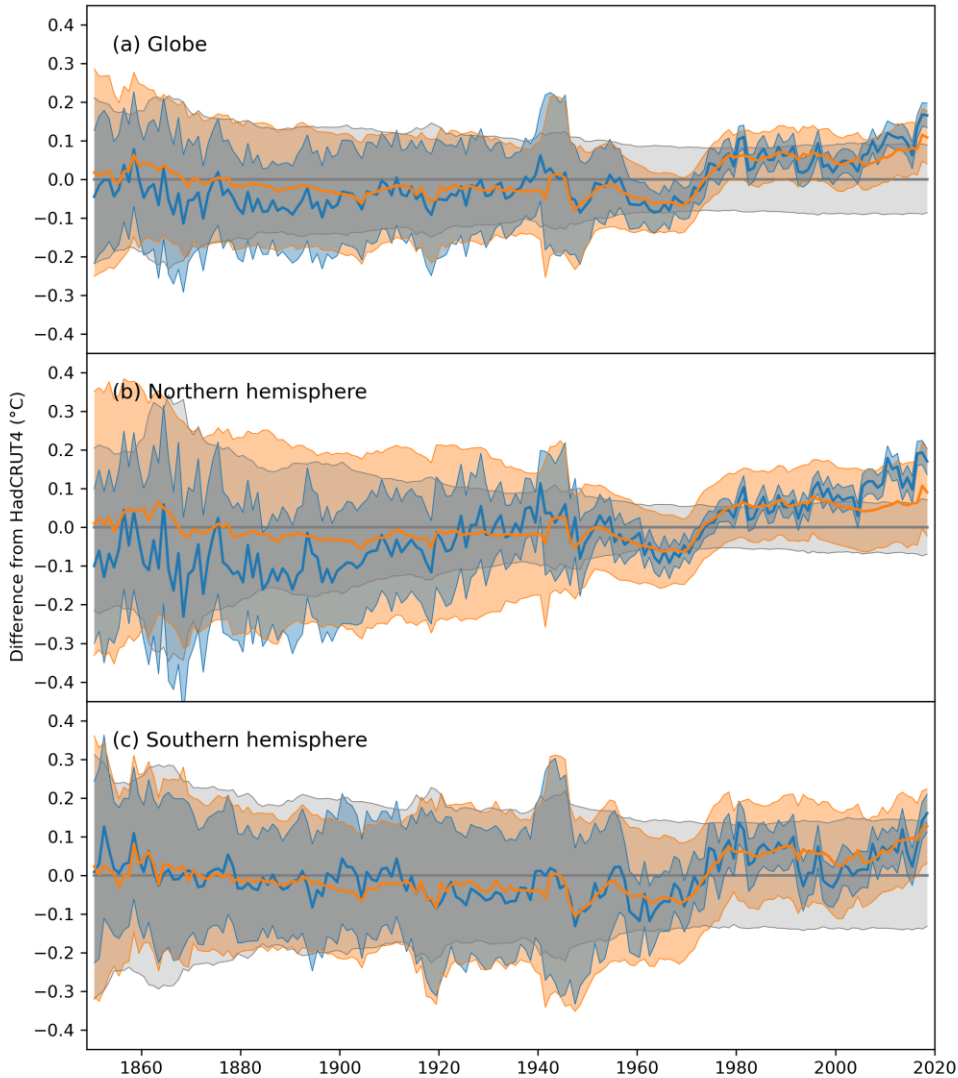
581

582

583 **4 Results**

584 4.1 Effects of updated data and methods in HadCRUT5

585



586

587 **Figure 1.** Annual average difference between HadCRUT.5.0.0.0 and HadCRUT.4.6.0.0 (°C),
 588 1850-2018. **(a)** Globe, **(b)** Northern Hemisphere and **(c)** Southern Hemisphere. Orange: non-
 589 infilled HadCRUT5. Blue: HadCRUT5 analysis. Solid lines: ensemble mean (HadCRUT.5.0.0.0)
 590 or median (HadCRUT.4.6.0.0). Orange/blue shading: 95% confidence interval determined by the
 591 ensemble spread and coverage uncertainty (the blue shading for the HadCRUT5 analysis lies
 592 mostly within the orange shading, where it appears as a darker grey due to the overlap). Light
 593 grey shading: 95% confidence interval on HadCRUT.4.6.0.0. Global means have been calculated
 594 by averaging hemispheric anomaly series for northern and southern hemispheres with equal
 595 weighting given to each hemisphere.

596

597 Differences in global and hemispheric mean time series between HadCRUT4 (version
598 HadCRUT.4.6.0.0) and the HadCRUT5 non-infilled data set and HadCRUT5 analysis are shown
599 in Figure 1. The differences between the non-infilled HadCRUT5 and HadCRUT4 primarily
600 arise from updates to the SST observational bias assessment in HadSST4. The updated bias
601 corrections result in slightly cooler anomalies globally and in each hemisphere from the 1880s to
602 1970s. Anomalies are warmer from the 1980s onwards.

603 The most obvious difference is the relative warming of HadCRUT5 between around 1970 and
604 1980. This arises from improved estimates of biases in measurements made in ship engine rooms
605 at that time. Engine room measurements were biased warm in the 1960s with the warm bias
606 dropping over time, first between 1970 and 1980 and then again between the early 2000s and
607 present. There are also changes around the Second World War, where changes to the
608 assumptions made in HadSST4 about how measurements were taken shifted the mean and
609 broadened the uncertainty range, reflecting the lack of knowledge of biases during this difficult
610 period (Kennedy et al., 2019).

611 Northern hemisphere uncertainty estimates for the non-infilled HadCRUT5 are slightly wider
612 than those of Morice et al. (2012). This results from a combination of the changes in the SST bias
613 adjustment model and the adoption of ERA5 as the reference data set for coverage uncertainty
614 calculations (Section 3.5). This change of reference data set typically gives wider uncertainty
615 estimates in the northern hemisphere for similar observational coverage. The reverse is true in
616 the southern hemisphere, with similar or slightly smaller coverage uncertainty estimates for the
617 non-infilled HadCRUT5. This reflects differences in regional variability in sparsely observed
618 regions between reanalysis products.

619 Further differences from HadCRUT4 can be seen in the HadCRUT5 analysis. Temperatures in
620 the latter decades of the 19th century are on average cooler than in the non-infilled HadCRUT5
621 data set in the global and northern hemisphere series. Temperatures in the 21st century are on
622 average warmer than those in the non-infilled HadCRUT5, primarily due to estimation of
623 additional areas of warm anomalies in high latitude regions in the northern hemisphere, including
624 use of air temperature anomalies over sea ice inferred from land stations. Rebalancing the
625 representation of land and marine regions also affects average temperatures throughout the
626 record. This is consistent with previous studies that adopt local interpolation methods (Cowtan &
627 Way, 2014; Karl et al., 2015; Lenssen et al., 2019). Together these features result in greater
628 warming throughout the 20th and 21st centuries in the HadCRUT5 analysis than is indicated by
629 the non-infilled data set. However, for any given year, the effect of the reconstruction may be to
630 produce either a warmer or cooler annual average and is dependent on variability in
631 reconstructed regions that were not well represented in HadCRUT4 (see also Figure 5 (b) and
632 (d)). Global and northern hemisphere HadCRUT5 analysis series fall outside the upper 95%
633 uncertainty limit of HadCRUT4 in the 21st century but rarely depart from the uncertainty range
634 of the HadCRUT5 non-infilled dataset, which includes the updated HadSST4 bias adjustments
635 and has wider northern hemisphere coverage uncertainty ranges.

636 The uncertainty range for the HadCRUT5 analysis is narrower than that for the non-infilled data
637 set, as the infilling effectively reduces the coverage uncertainty by filling gaps in the data and
638 accounting for the non-uniform distribution of observations. The effect of this can be clearly
639 seen in the Southern Hemisphere (Figure 1) where the narrowing of the uncertainty range before
640 the 1950s is much less than after the 1950s, when routine monitoring on the Antarctic continent
641 started, and coverage of the HadCRUT5 analysis thereafter approaches 100%.

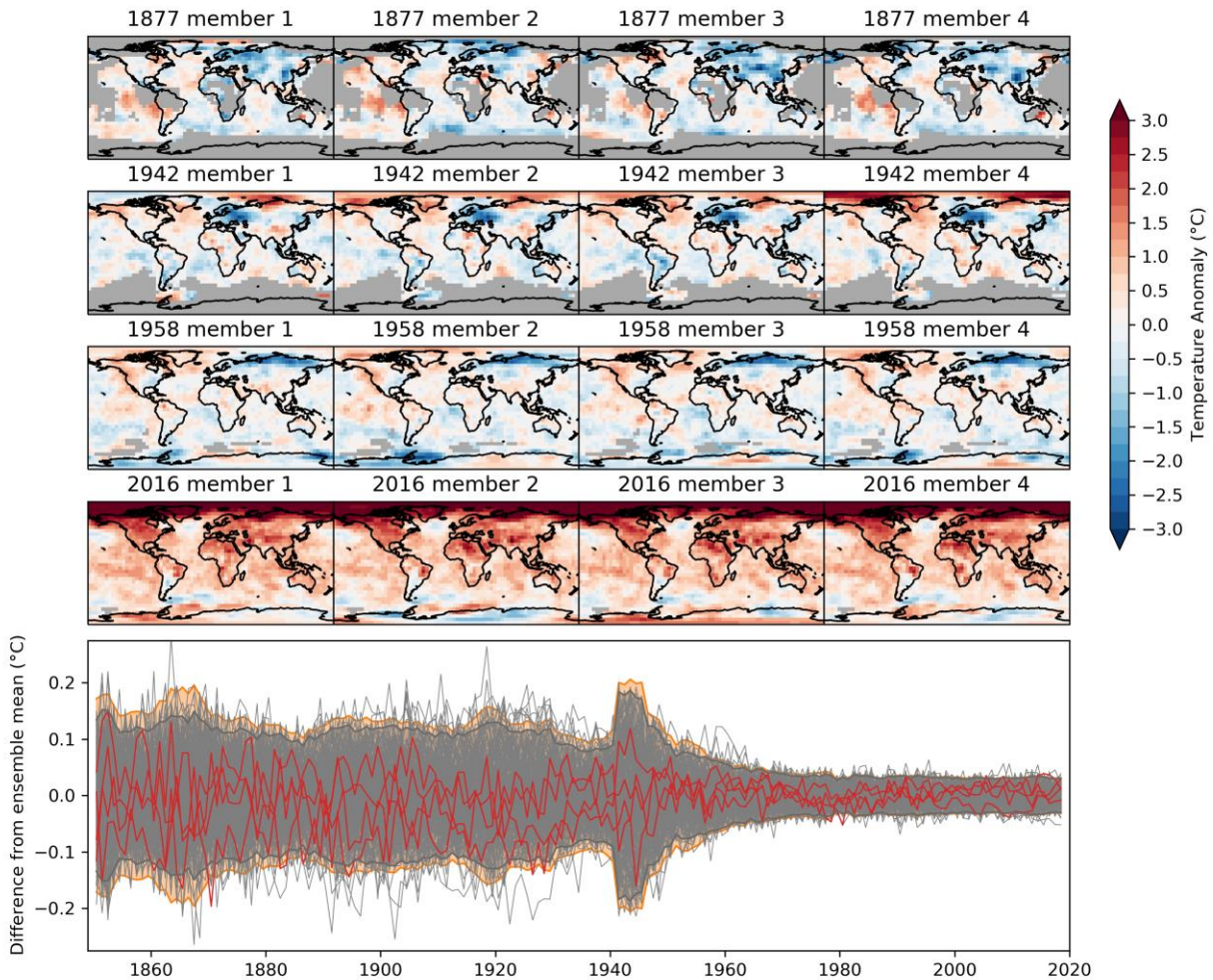
642 As discussed in Section 3.1, the error model structure for the non-infilled HadCRUT5 data set is
643 the same as in Morice et al. (2012), with observational bias adjustment uncertainties encoded
644 into the ensemble and separate measurement and sampling uncertainty information provided and
645 propagated into the uncertainty ranges on the hemispheric and global averages shown in Figure
646 1. The approach adopted for the HadCRUT5 analysis differs in including the effects of
647 measurement and sampling uncertainties in the ensemble while also sampling from the spatial
648 analysis uncertainty. Examples of HadCRUT5 analysis ensemble members are shown in Figure
649 2.

650 There is little change in the HadCRUT5 analysis ensemble spread for global or hemispheric
651 averages from the 1970s onwards, reflecting the spread in the underlying SST ensemble and the
652 relatively stable spatial sampling during this period. The ensemble spread in the global average
653 in the 1940s is similar to that prior to the 1870s, though in the 1940s, this spread arises
654 predominantly from uncertainty in the SST biases, whereas prior to the 1870s, the spread is
655 largely due to uncertainty in the spatial field estimates due to limited observational sampling of
656 the globe.

657 There is coherent spatial structure in the deviations of ensemble member fields from the
658 ensemble mean. This results from uncertainty in the spatial analysis and its estimation from
659 uncertain observations. Some ensemble members may be cool while others are warm in regions
660 where uncertainty is high (for example see differences between ensemble members in Antarctica
661 in Figure 2). The additional coverage uncertainty arising from masked regions is a relatively
662 smaller component of the total uncertainty as a result of the increased coverage in the
663 HadCRUT5 analysis fields and the inclusion of reconstruction uncertainty within the ensemble.
664 On multi-annual timescales, the uncertainty in observational bias adjustments becomes
665 prominent. This is reflected in persistently warm or cool departures from the ensemble mean in
666 global and regional diagnostics over many years for individual ensemble members (for example
667 see ensemble series in Figure 2).

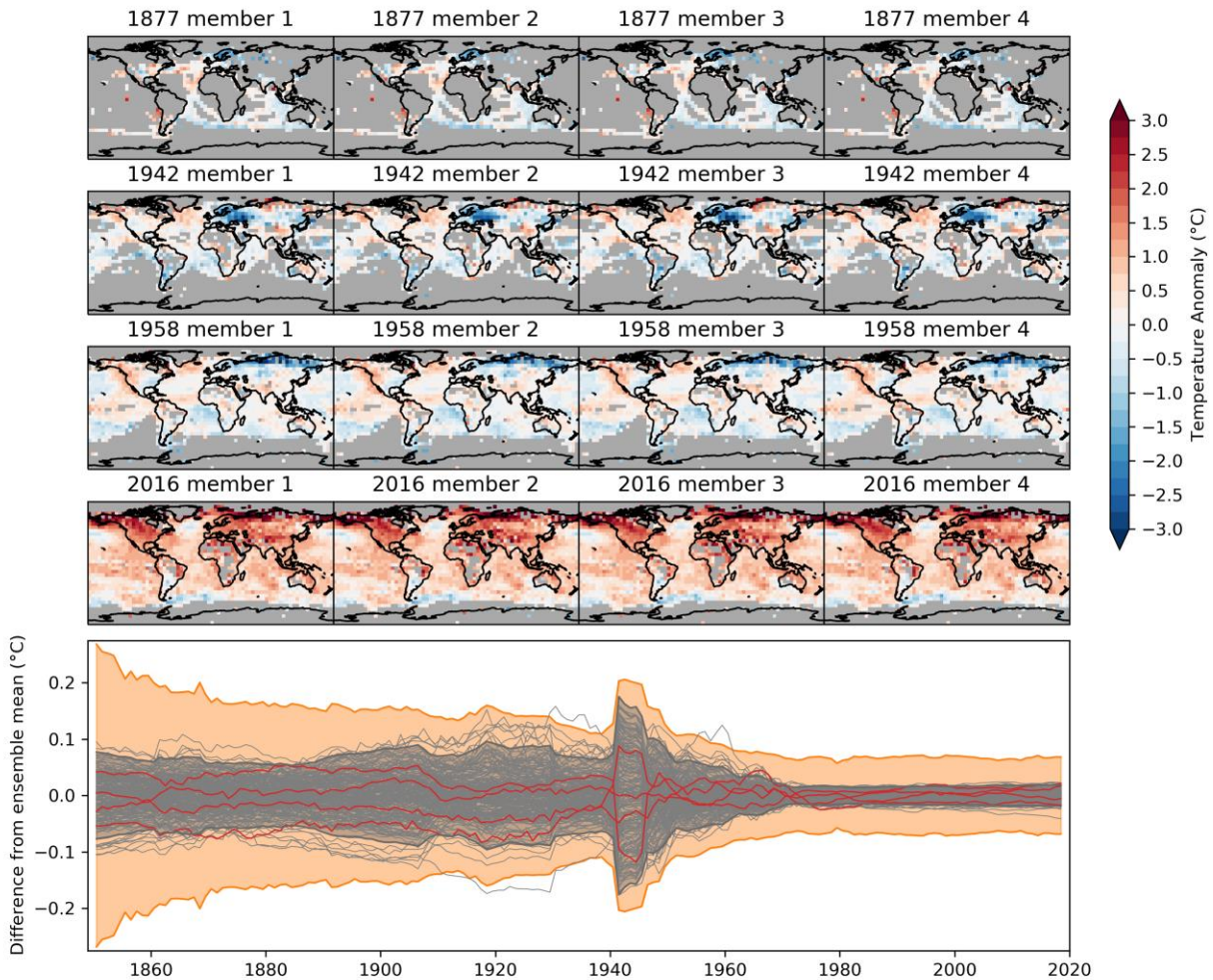
668 Non-infilled HadCRUT5 ensemble members are shown in Figure 3, matching those shown for
669 the HadCRUT5 analysis in Figure 2. HadCRUT5 analysis fields have greater spatial extent than
670 the non-infilled dataset and are also smoother as a result of measurement and sampling
671 uncertainties being taken into account within the analysis framework. In regions of few, scattered
672 observations, infilled analysis fields have much greater extent but also show diversity in
673 reconstructed anomaly patterns, reflecting uncertainty in the reconstruction in these sparsely
674 observed regions.

675 Uncertainty ranges for the global average temperature series in Figure 3 show the ensemble
676 spread in relation to the full uncertainty range, accounting for all quantified sources of
677 uncertainty. While the HadCRUT5 analysis and non-infilled data set quantify uncertainty from
678 the same error sources, the HadCRUT5 analysis encodes a greater portion of the uncertainty into
679 the ensemble, whereas the non-infilled ensemble only samples uncertainties that are most
680 important over multi-decadal time scales. The ensemble for the non-infilled HadCRUT5 data set
681 samples the uncertainty associated with observational bias adjustments, with structure that is
682 relevant to multi-decadal climate assessments. Unlike the HadCRUT5 analysis, measurement
683 and sampling uncertainties that are relevant at shorter time scales are not encoded into the
684 ensemble and are instead provided as auxiliary information. Uncertainty from incomplete global
685 coverage of the observing network is a greater portion of the total uncertainty for the non-infilled
686 data set. In contrast, for the HadCRUT5 analysis, the uncertainty from incomplete global
687 coverage is divided between the analysis ensemble spread in reconstructed regions and a smaller
688 coverage uncertainty term relating to regions that are masked.



689

690 **Figure 2.** HadCRUT5 analysis ensemble members. Upper panel: annual average temperature
 691 anomaly ($^{\circ}\text{C}$, relative to 1961-90) for 1877, 1942, 1958 and 2016 in four example ensemble
 692 members. Lower panel: ensemble spread in global mean ($^{\circ}\text{C}$), 1850-2018. The difference
 693 between each ensemble member and the ensemble mean is shown by the grey lines, with the first
 694 four ensemble members (corresponding to the maps above) highlighted in red. Grey shading:
 695 95% confidence interval determined by the ensemble spread. Orange: full uncertainty range
 696 adding the additional coverage uncertainty term. Global means have been calculated by
 697 averaging anomalies for northern and southern hemispheres for each ensemble member. Maps
 698 require six months of data within a year for a grid cell average to be plotted.



699

700 **Figure 3.** As Figure 2, but for the HadCRUT5 non-infilled dataset. Upper panel: annual average
 701 temperature anomaly ($^{\circ}\text{C}$, relative to 1961-90) for 1877, 1942, 1958 and 2016 in four example
 702 ensemble members. Lower panel: ensemble spread in global mean ($^{\circ}\text{C}$), 1850-2018. The
 703 difference between each ensemble member and the ensemble mean is shown by the grey lines,
 704 with the first four ensemble members (corresponding to the maps above) highlighted in red. Grey
 705 shading: 95% confidence interval determined by the non-infilled ensemble spread. Orange: full
 706 uncertainty range including additional measurement and sampling uncertainty terms, that are not
 707 sampled by the non-infilled ensemble, and the coverage uncertainty term. Global means have
 708 been calculated by averaging anomalies for northern and southern hemispheres for each
 709 ensemble member. Maps require six months of data within a year for a grid cell average to be
 710 plotted.

711

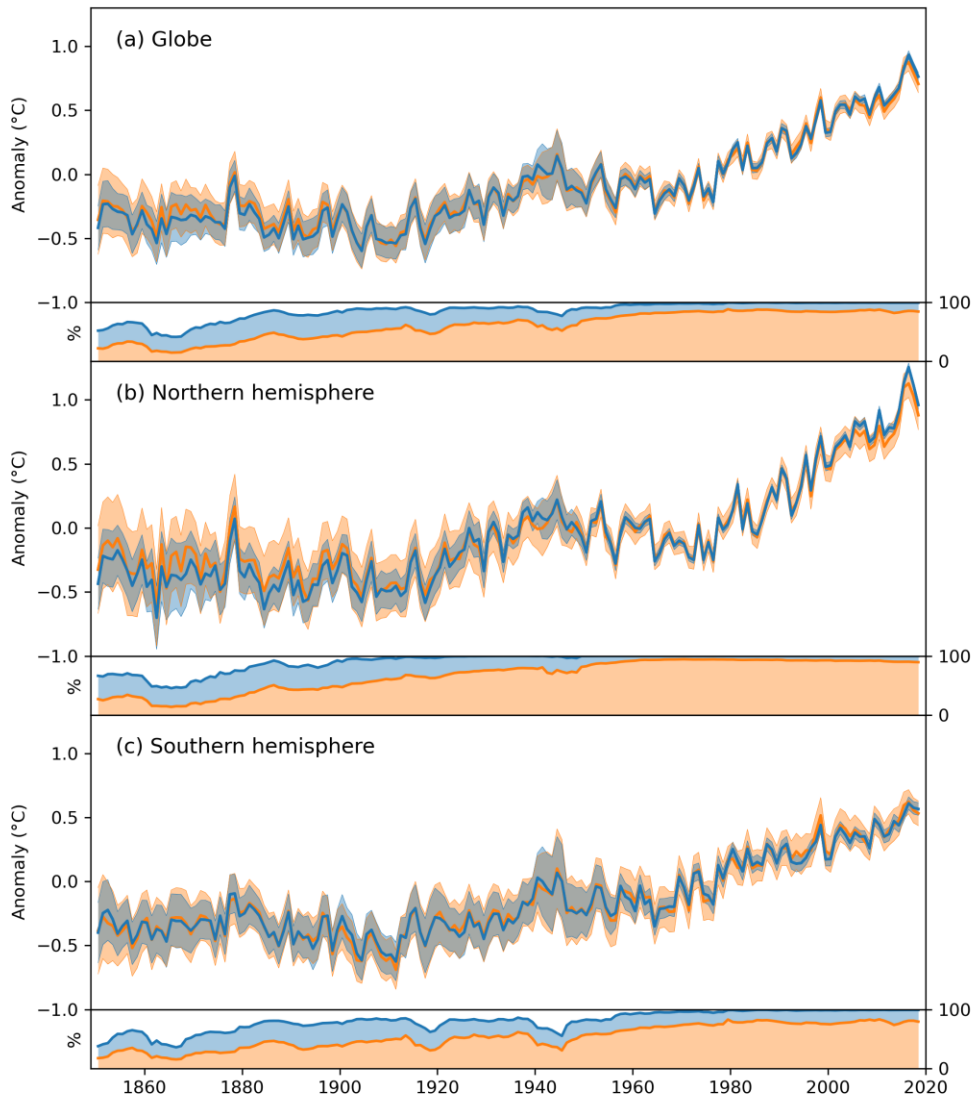
712 4.2 Global, hemispheric and regional series

713 Annual global and hemispheric average temperature anomaly series for HadCRUT5 are shown in
714 Figure 4, along with the fraction of regional data coverage represented in the non-infilled dataset
715 and the HadCRUT5 analysis.

716 Areal data coverage in the HadCRUT5 analysis grids first reaches 90% in the 1900s, with two
717 subsequent drops in coverage in the late 1910s and early 1940s associated with the two world
718 wars. Northern hemisphere coverage exceeds 99% in the early 1920s and reaches 100% in the
719 mid-1950s. Uncertainty in southern hemisphere temperatures is greatest in the period prior to the
720 establishment of a sustained Antarctic monitoring network in the 1950s (see also Figure 5 (a)),
721 after which global coverage exceeds 97% in the 1960s. The spatial extent of the observing
722 network in the southern hemisphere is also a prominent contribution to uncertainty in global
723 average series prior to the 1950s. Global coverage of the analysis fields is typically not complete
724 even in modern years due to an absence of sustained observation in the southern South Pacific,
725 and the nearby Southern Ocean and Antarctic.

726 Southern Hemisphere anomalies are cooler in the HadCRUT5 analysis in the 1990s from around
727 1992, particularly in 30-60S (Figure 5 (b)). The observing network is less dense in these regions,
728 with regular shipping covering only the equatorward half of the latitude band, leading to
729 differences between non-infilled HadCRUT5 and the HadCRUT5 analysis. Variability in the
730 regional time series (Figure 5) is smaller in the early record in the HadCRUT5 analysis than the
731 non-infilled dataset, particularly in the high latitude regions as a result of reduced uncertainty
732 from spatial sampling in the HadCRUT5 analysis.

733

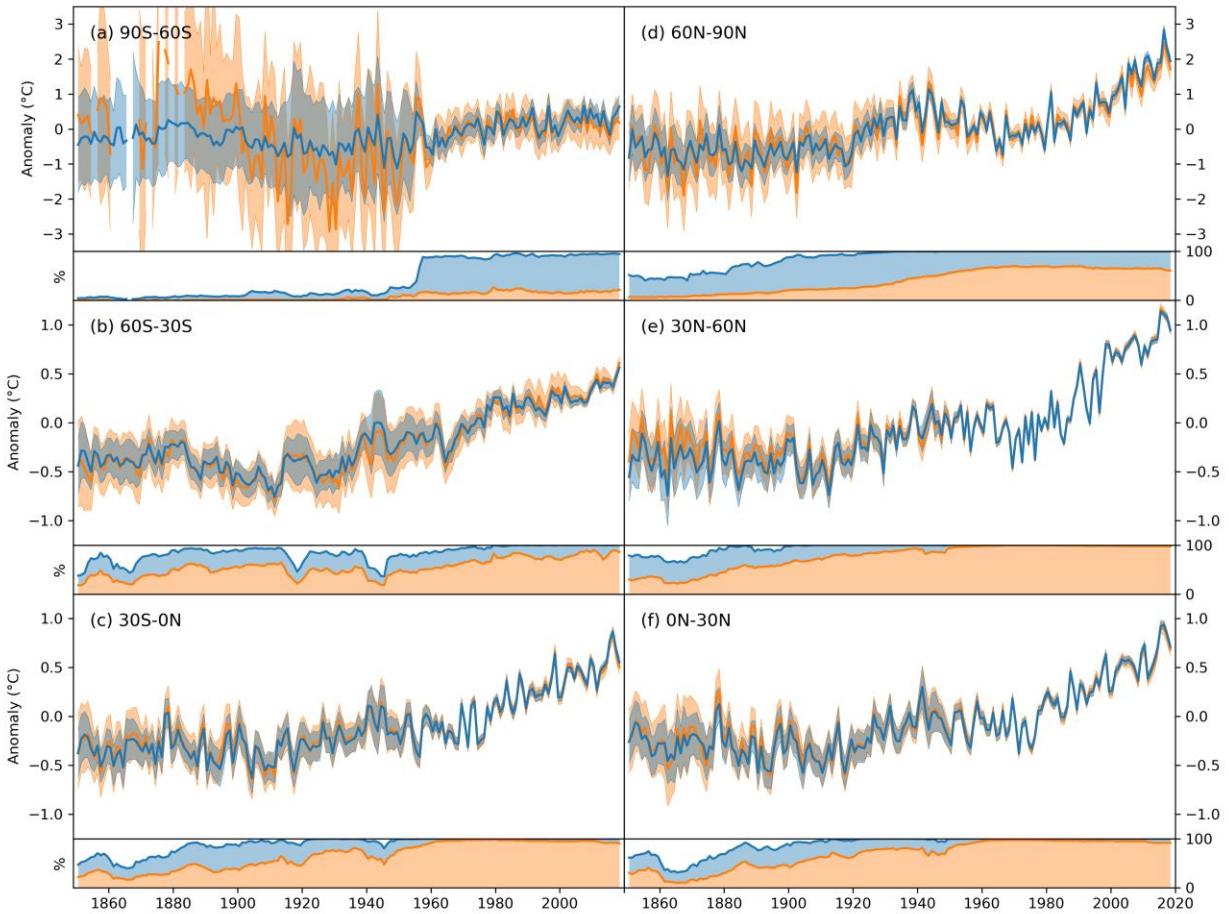


734

735 **Figure 4.** Comparison between the HadCRUT5 analysis and non-infilled data set. **(a)** Globe, **(b)**
 736 Northern Hemisphere and **(c)** Southern Hemisphere. Upper panel in each pair: annual average
 737 temperature anomaly ($^{\circ}\text{C}$, relative to 1961-90), 1850-2018. Lower panel in each pair: percentage
 738 of area covered by data in each annual average. Orange: non-infilled HadCRUT5 data set. Blue:
 739 HadCRUT5 analysis. Solid lines: ensemble mean. Orange/blue shading: 95% confidence
 740 interval. Global means have been calculated by averaging anomalies for northern and southern
 741 hemispheres.

742

743



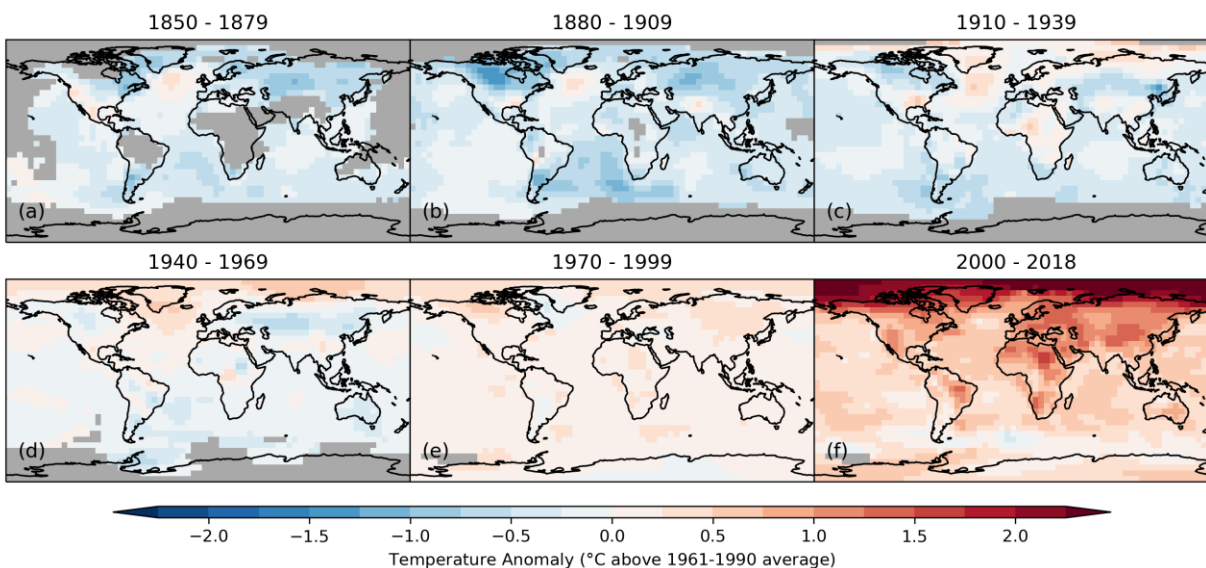
744

745 **Figure 5.** Comparison between the HadCRUT5 analysis and non-infilled data set. **(a)** 90°S-60°S,
 746 **(b)** 60°S-30°S, **(c)** 30°S-0°N, **(d)** 60°N-90°N, **(e)** 30°N-60°N and **(f)** 0°N-30°N. Upper panel in
 747 each pair: annual average temperature anomaly (°C, relative to 1961-90), 1850-2018. Lower
 748 panel in each pair: percentage of area covered by data in each annual average. Orange: non-
 749 infilled HadCRUT5 data set. Blue: HadCRUT5 analysis. Solid lines: ensemble mean.
 750 Orange/blue shading: 95% confidence interval.

751

752

753



754

755 **Figure 6.** Long-term average temperature anomaly ($^{\circ}\text{C}$, relative to 1961-90). (a) 1850-1879, (b)
 756 1880-1909, (c) 1910-1939, (d) 1940-1969, (e) 1970-1999 and (f) 2000-2018. Averages require at
 757 least one month per quarter, three quarters per year, and 50% of years per multi-year period.

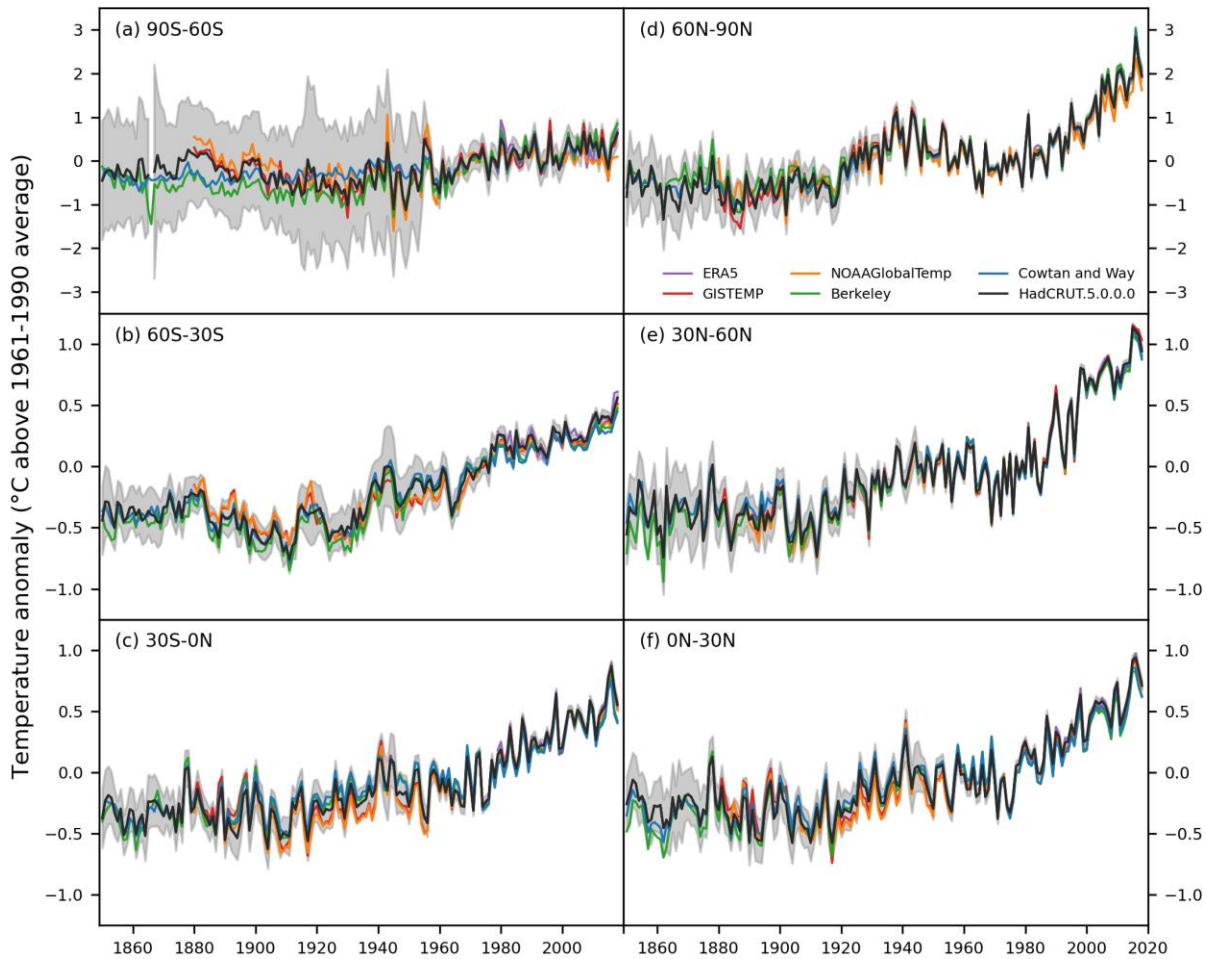
758 In regions where data are sparse, and hence uncertainty in surface temperature analyses is
 759 largest, data that might be used to validate the analyses is also highly limited. Here we have used
 760 the ratio of posterior to prior variances to remove regions with weak observational constraint (see
 761 Appendix for details). Despite restricting the reconstruction to regions that are locally
 762 constrained, there is a marked increase in the area of the globe represented by the HadCRUT5
 763 analysis in comparison to the non-infilled data set (see coverage timeseries in Figure 5 and
 764 example monthly fields in Figures S10 to S13 of the Supporting Information).

765 Figure 6 reveals the patterns of change in successive 30-year periods and the most recent 19
 766 years of the HadCRUT5 analysis. Even in these longer-term averages, there are regions that are
 767 particularly warm or cool relative to the global mean. The final panel for 2000-2018 illustrates
 768 the greater warming at high northern latitudes and over the land compared to the ocean. The
 769 surface waters of the Southern Ocean, in contrast, have warmed more slowly than many other
 770 areas. We also see one area of long-term cooling, to the south of Greenland and Iceland (Parker
 771 et al., 1994). 1880-1909 was a particularly cool period, with centers of low average anomalies in
 772 the South Atlantic, Canada and central Russia.

773

774 4.3 Comparisons with other analyses

775 Average temperature changes over the whole period of record in 30° latitude bands for a range of
 776 analyses are shown in Figure 7. These analyses include NOAA GlobalTemp v5 (Huang et al.
 777 2019), NASA GISTEMP v4 (Hansen et al., 2010; Lenssen et al. 2019), the Cowtan & Way
 778 analysis (Cowtan & Way, 2014), and the Berkeley Earth analysis (Rohde & Hausfather, 2020).
 779 The HadCRUT.5.0.0.0 analysis is also shown.



780

781 **Figure 7.** Comparison between long-term near-surface temperature data sets. Annual average
 782 temperature anomaly (°C, relative to 1961-90), 1850-2018. **(a)** 90°S-60°S, **(b)** 60°S-30°S, **(c)**
 783 30°S-0°N, **(d)** 60°N-90°N, **(e)** 30°N-60°N and **(f)** 0°N-30°N. Black: HadCRUT5 analysis
 784 ensemble mean. Pink: ERA5. Red: GISTEMP. Orange: NOAA GlobalTemp. Green: Berkeley
 785 Earth. Blue: Cowtan & Way. Grey shading: 95% confidence interval on the HadCRUT.5.0.0.0
 786 analysis determined by the ensemble spread and coverage uncertainty.

787 All of the analyses shown use spatial infilling. Cowtan & Way and Berkeley Earth use
 788 interpolation methods based on a statistical model of local covariance structure (although within
 789 a more complex statistical model of global temperature variation in the Berkeley Earth analysis).
 790 NOAA GlobalTemp uses a model of spatially-varying local patterns of temperature variability.

791 GISTEMP employs a distance-weighted interpolation for land based meteorological station data
792 and uses the same large-scale analysis of sea-surface temperatures used in NOAAGlobalTemp.
793 GISTEMP, Cowtan & Way and Berkeley are each close to globally complete since the 1950s
794 while the NOAAGlobalTemp data set does not extend into data-sparse polar regions.

795 The analyses are most similar in regions with the densest observational coverage, such as in the
796 northern mid-latitudes (Figure 7 (e)). Where observational coverage is lowest, the analyses
797 become sensitive to assumptions underpinning reconstruction methods. For example,
798 NOAAGlobalTemp reconstructs fields through low-frequency smoothing and a model of
799 dominant spatial patterns of variability, while methods based on local covariance structure may
800 tend toward a field mean in the case of Cowtan & Way, Berkeley, or the HadCRUT5 analysis
801 ensemble mean, or towards the anomalies observed at nearby locations for the GISTEMP land
802 analysis method. The analyses also differ in how regions that are distant from observed locations
803 are included or are masked.

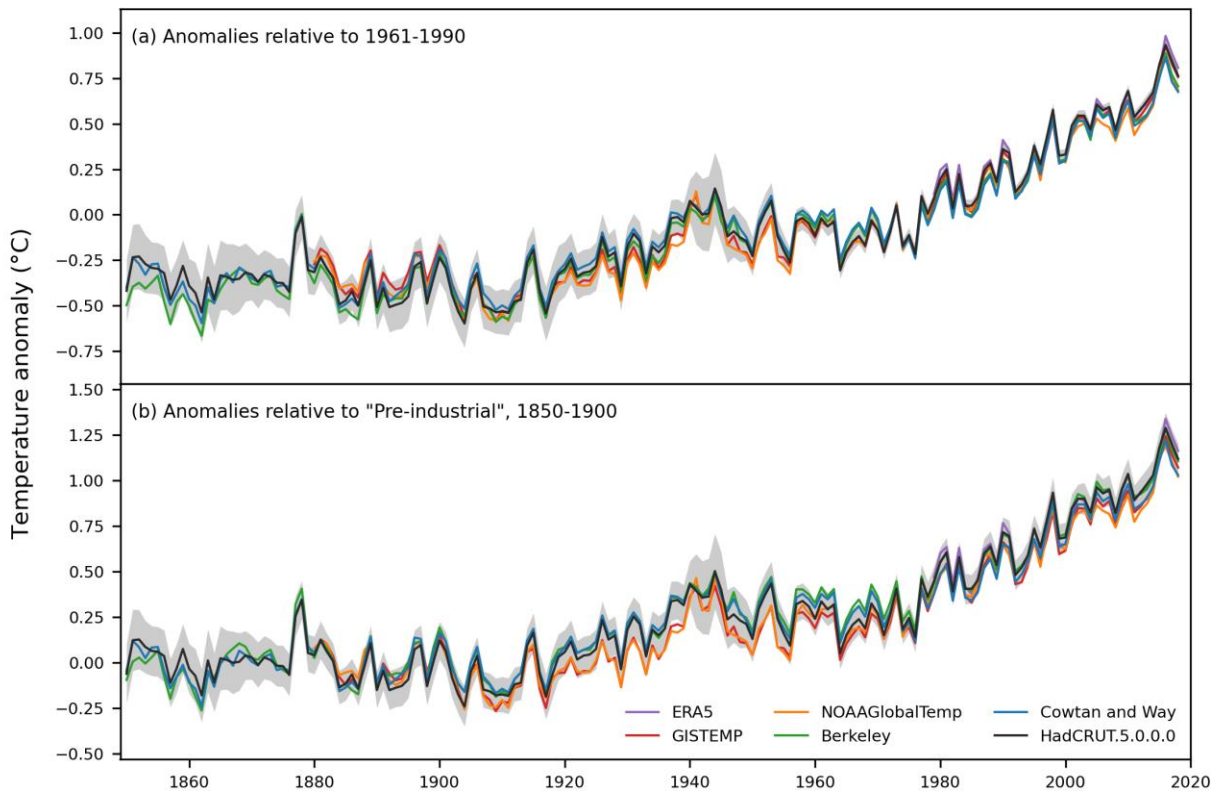
804 The HadCRUT5 analysis method is closely related to the method used in Cowtan & Way but
805 differs in three key aspects. First, it accounts for the spatial variation in data uncertainty as well
806 as the estimated measurement and sampling error covariances. This is particularly important for
807 the oceans, where less-reliable ship data are combined with more accurate data from drifting and
808 moored buoys. Second, the spatial analysis method is used to make improved temperature
809 estimates at all locations, not just grid cells without data. Third, by using a full covariance model
810 for both the temperature field and the observational uncertainty within a Bayesian analysis
811 framework, it is possible to sample from the posterior of the distribution to generate a consistent
812 ensemble data set that combines all known sources of uncertainty whilst respecting the estimated
813 covariance structure of the temperature anomaly field.

814 The differences between the HadCRUT5 analysis ensemble mean and Cowtan & Way in the post
815 1950 period, are largely due to changes in the estimated SST biases. As Berkeley Earth shows
816 similar differences and uses the same SST data set as Cowtan & Way, we can infer that changes
817 in the estimated SST biases are the key difference here as well. The changes in SST bias
818 estimates are larger in the more sparsely observed regions – the tropics and southern hemisphere
819 – where there are fewer ships, so changes in assumptions about observing practice of a few
820 countries can have a proportionately larger effect.

821 Differences between HadCRUT5 and the ERSST-based data sets, GISTEMP and
822 NOAAGlobalTemp are also largely due to differences in estimated SST biases. In particular,
823 ERSST tends to be cooler than HadSST4 from the early 20th century to the start of the Second
824 World War and from the end of the war to around 1955; this difference is associated with
825 uncertainty in the estimated biases associated with bucket measurements, particularly in the
826 Southern Hemisphere and the tropics. From the 1960s, agreement between HadSST4 and ERSST
827 is better, though there is a notable cooling of ERSST relative to HadSST4 in the early 1990s
828 associated with a relative cooling of marine air temperature compared to SST (see Kennedy et al.
829 2019). From the late 1990s onwards, both ERSSTv5 and HadSST4 show good relative stability
830 compared to instrumentally homogeneous data sets (Hausfather et al., 2017; Kennedy et al.,
831 2019). Notable structural uncertainty remains in early SST records.

832 Differences can be seen in the first half of the 20th century between
 833 GISTEMP/NOAAGlobalTemp and Cowtan & Way/HadCRUT5 over the latitude band 0°N-30°S
 834 with GISTEMP/NOAAGlobalTemp cooler (Figure 7 (c)). Regional differences over land partly
 835 result from differences in homogenization and the underlying station data sets. HadCRUT5 uses
 836 homogenized station data (from CRUTEM5), as provided by national meteorological services or
 837 research projects. Other datasets include automated homogenization algorithms (Huang et al.,
 838 2019; Menne et al., 2018; Rohde et al., 2013b). This may result in regional differences between
 839 data sets, particularly where the measurement network is less dense and, as a consequence, there
 840 is greater uncertainty in homogenization.

841



842

843 **Figure 8.** Comparison of annual global average temperature anomaly series (°C) relative to two
 844 baselines: (a) 1961-1990 and (b) 1850-1900, taken as representative of pre-industrial conditions.
 845 Black: HadCRUT5 analysis ensemble mean. Pink: ERA5. Red: GISTEMP. Orange:
 846 NOAAGlobalTemp. Green: Berkeley Earth. Blue: Cowtan and Way. Grey shading: 95%
 847 confidence interval on the HadCRUT5 analysis determined by the ensemble spread only. Global
 848 means have been calculated for each data set by averaging anomalies for northern and southern
 849 hemispheres. For all datasets except for ERA5, anomaly series are computed by adjusting
 850 monthly time series to the appropriate baseline using data available in the anomaly reference
 851 period before averaging to annual series. ERA5 timeseries are shifted to match the 1981-2010
 852 average for the HadCRUT5 analysis series, due to insufficient data in the climatology periods to

853 compute anomalies. Anomaly series and uncertainties provided by the dataset producers using
854 each dataset's native methods are shown in Supporting Information Figure S9.

855 Temperature changes relative to the average over the late 19th century are shown in Figure 8.
856 The 51-year period 1850-1900 is often considered for practical purposes to be representative of
857 pre-industrial conditions. This approximation of pre-industrial temperatures is consistent with
858 that adopted in IPCC AR5 (Hartmann et al., 2013) and IPCC SR1.5 (Allen et al., 2018), noting
859 that any choice of period is a compromise, with natural variability and forcing playing a role
860 (Hawkins et al., 2017). For analyses that do not extend back to 1850 (NOAAGlobalTemp and
861 GISTEMP), 1880 to 1900 is used as the reference period here. By referencing the time series to
862 this early period, the spread of temperature anomalies later in the series is increased. This
863 increased spread reflects uncertainty in temperatures in the early reference period and not
864 uncertainty in recent temperature changes. On the global mean, the analyses are remarkably
865 consistent with one another despite the differences in their construction.

866 **5 Conclusions**

867 An updated data set of global near-surface temperature change, HadCRUT5, is presented.
868 Updates in the CRUTEM5 dataset have expanded the underlying land station series and provided
869 additional data quality checks. Updates in HadSST4 have brought improved understanding of the
870 evolution of the marine observing network, contributing improved bias adjustments and
871 uncertainty estimates. These are combined both in a non-infilled data set and in a new ensemble
872 statistical analysis that provides a more spatially complete assessment of global and regional
873 changes and uncertainty therein.

874 The new HadCRUT5 analysis ensemble samples a greater range of the quantified uncertainties
875 than our previous assessment (Morice et al., 2012). Uncertainties arising from systematic errors
876 associated with observational methods, measurement and sampling errors and spatial analysis
877 uncertainty are all encoded into the expanded 200-member ensemble, communicating the major
878 known sources of uncertainty in an easily accessible way.

879 Time series of globally averaged temperature anomalies show greater 21st century warming for
880 the HadCRUT5 analysis than for the HadCRUT5 non-infilled data set. The increased warming is
881 predominantly associated with improved representation of the rapidly warming but sparsely
882 observed high latitudes of the northern hemisphere. This finding is consistent with other
883 independently-produced statistical analyses of global temperature changes and is also consistent
884 with temperature changes observed in reanalysis data sets that assimilate observational data into
885 a numerical weather prediction model (Kobayashi et al., 2015; Gelaro et al., 2017; Blunden &
886 Arndt, 2019; Hersbach et al., 2020).

887 The HadCRUT5 analysis indicates that globally averaged temperatures in the second half of the
888 19th century were on average cooler than estimates based on non-infilled HadCRUT5. This is
889 also consistent with assessments based on other independently produced statistically infilled
890 analyses. Combined with the evidence of increased warming in recent years, the infilled analyses
891 indicate that warming since the 19th century is likely greater than is indicated by HadCRUT4 as
892 a result both of observational sampling in the non-infilled data set and of updates to our

893 understanding of biases in sea-surface temperature measurements resulting from changes in the
894 make-up of the marine observing network.

895 There is, however, uncertainty in our understanding of 19th century temperatures resulting from
896 limitations in observational sampling, particularly in the southern hemisphere, and uncertainty
897 associated with residual observational biases. Uncertainty remains in the early instrumental
898 record in locations for which observational data are not available to inform the analysis. This is
899 most evident in the Antarctic, the Arctic and regions of the southern hemisphere land, prior to the
900 establishment of permanent observing sites.

901 Methodological choices in representation of data sparse regions in different data sets lead to
902 differences between global and regional average temperature time series. The impacts of these
903 choices are most evident in regions and at times in which the observational data required to
904 constrain the analysis is limited or unavailable, particularly in regions of the southern hemisphere
905 in the early record. The spread of 19th century temperature analyses produced by different
906 monitoring centers in part reflects the sensitivity to differences in methods used. These methods
907 assume different statistical models for the data; therefore, the differences between analyses are
908 not necessarily captured by the uncertainty estimates of any single method.

909 The updated analysis methods assist in mitigation of the impacts of low availability of
910 observational data in data sparse regions. We anticipate that an extension, in potential future
911 work, of the analysis covariance model to describe regional variation in variability would further
912 improve the analysis temperature fields and uncertainty estimates. However, digitization of as
913 yet unavailable observations and submission of these to open archives continues to be invaluable
914 to improve regional data coverage and reduce uncertainty further.

915 The use of marine air temperature observations has recently been proposed to reconcile
916 differences between datasets produced as a blend of SST and air temperature observations and
917 model-based studies using near-surface air temperatures over ocean (Cowtan et al., 2015;
918 Richardson et al., 2016). However, uncertainties in observed long-term changes in marine air
919 temperature and their differences from observed SSTs are important to understand (Kennedy et
920 al. 2019, Chan and Huybers 2019, Chan et al. 2019), and the marine air temperature observing
921 network is less robust than that for SST and is in long-term decline (Berry & Kent, 2017).
922 Challenges also remain in monitoring near-surface temperature changes in the cryosphere, given
923 sparse observational coverage and changes in sea-ice extent, with impacts on downstream
924 assessments (Richardson et al., 2018).

925 Relative biases in sea-surface temperature measurements arise from differences in measurement
926 methods and instrumentation. Such biases change regionally and over time with gradual as well
927 as abrupt changes in the composition of the observing network or underlying databases. The
928 characteristics of different bias adjustment schemes can be seen in the differences between
929 analyses, broadly grouping data sets into those (GISTEMP, Lenssen et al. (2019) and
930 NOAAGlobalTemp, Huang et al. (2019)) that adopt the ERSST v5 dataset (Huang et al., 2017),
931 those (Cowtan & Way (2014) and Berkeley Earth (Rohde & Hausfather, 2020)) that adopt
932 HadSST3 (Kennedy et al., 2011a and b), and that which uses the improved HadSST4 data set
933 (Kennedy et al., 2019), as is documented here. Differences between bias adjustments applied in
934 each data set are smaller than the assessed adjustments themselves, which result in a net

935 reduction in observed warming compared to unadjusted measurements (Kennedy et al., 2019).
936 Nevertheless, differences in SST bias assessments feature prominently as a source of difference
937 between studies and remain a key uncertainty in assessing long-term change (Kent et al., 2017).

938 Despite methodological differences, temperature series derived from different analyses are in
939 good agreement, generally lying within the assessed uncertainty range of the HadCRUT5
940 analysis. Updates in HadCRUT5 bring our estimates of global and hemispheric series closer to
941 those of other recent studies. Remaining differences between estimates are understood to
942 predominantly arise from differences in spatial analysis methods applied and differences in how
943 each analysis accounts for changes in marine observing methods.

944 **Acknowledgments**

945 CPM, JJK, NAR, JPW, EH, REK, RJHD and IRS were supported by the Met Office Hadley
946 Centre Climate Programme funded by BEIS and Defra. TJO and PDJ were supported by UK
947 NERC (grant number NE/N006348/1, SMURPHS). The CRUTEM5 dataset was developed as
948 part of the UK National Centre for Atmospheric Science. We thank three anonymous reviewers
949 for their constructive comments which improved the manuscript.

950

951 **Data access**

952

953 The gridded temperature anomalies, the global and hemispheric timeseries and their uncertainty
954 intervals will be available from the Met Office website (<https://www.metoffice.gov.uk/hadobs>).
955 HadCRUT5 data will be archived for long term preservation and reuse as part of the HadCRUT
956 catalogue at CEDA <https://catalogue.ceda.ac.uk/uuid/f7189fabb084452c9818ba41e59ccabd>. The
957 CEDA archive of the HadCRUT.5.0.0.0 data can be accessed from
958 <https://catalogue.ceda.ac.uk/uuid/b9698c5ecf754b1d981728c37d3a9f02>.

959

960 ERA5 was obtained from the Copernicus Climate Change Service (C3S) (2017): ERA5: Fifth
961 generation of ECMWF atmospheric reanalyses of the global climate. Copernicus Climate
962 Change Service Climate Data Store (CDS), date of access: 28/11/2019,
963 <https://cds.climate.copernicus.eu>.

964

965 HadISST.2.2.0.0 was accessed on 11/12/2019 from
966 <https://www.metoffice.gov.uk/hadobs/hadisst2/>.

967

968 The HadSST.4.0.0.0 ensemble is available from <https://www.metoffice.gov.uk/hadobs/hadsst4/>.

969

970 CRUTEM5 data will be available from <https://www.metoffice.gov.uk/hadobs> and the CRUTEM
971 collection at CEDA <https://catalogue.ceda.ac.uk/uuid/eeabb5e1ff2140f48e76ea1ffda6bb48>. The
972 CEDA archive of the CRUTEM.5.0.0.0 data can be accessed from
973 <https://catalogue.ceda.ac.uk/uuid/901f576daca4e049630ab879d6fb476>.

974

975 HadCRUT.4.6.0.0 is available from <https://www.metoffice.gov.uk/hadobs/hadcrut4/>.

976

977 GISTEMP version 4 was accessed on 17/11/2019 at 15:45 GMT from
978 <https://data.giss.nasa.gov/gistemp/>.

979
 980 NOAAGlobalTemp version 5 was accessed on 15/10/2019 at 07:07 GMT from
 981 [https://www.ncdc.noaa.gov/noaa-merged-land-ocean-global-surface-temperature-analysis-](https://www.ncdc.noaa.gov/noaa-merged-land-ocean-global-surface-temperature-analysis-noaaglobaltemp-v5)
 982 [noaaglobaltemp-v5](https://www.ncdc.noaa.gov/noaa-merged-land-ocean-global-surface-temperature-analysis-noaaglobaltemp-v5).

983
 984 Berkeley Earth was accessed on 17/11/2019 at 16:25 GMT from [https://berkeleyearth.org/data-](https://berkeleyearth.org/data-new/)
 985 [new/](https://berkeleyearth.org/data-new/).

986
 987 Cowtan and Way was accessed on 14/10/2019 at 10:40 GMT from [https://www-](https://www-users.york.ac.uk/~kdc3/papers/coverage2013/series.html)
 988 [users.york.ac.uk/~kdc3/papers/coverage2013/series.html](https://www-users.york.ac.uk/~kdc3/papers/coverage2013/series.html).

989

990 **Appendix A: Details of spatial analysis methods**

991 *A.1 Modelling the temperature anomaly field as a Gaussian process*

992 Here we describe the methods used to construct the HadCRUT5 analysis. The method described
 993 in this section follows the Gaussian process method with explicit basis functions, described in
 994 Rasmussen & Williams (2006). The methods for analysis hyperparameter estimation are
 995 described in Appendix A.2. Appendix A.3 describes application to the non-infilled land air
 996 temperature and sea surface temperature ensemble grids, including methods for sampling
 997 analysis uncertainties. Regional masking of the analyses is described in Appendix A.4.

998 For a monthly temperature anomaly field \mathbf{g} , we model a vector of gridded temperature anomaly
 999 observations \mathbf{y} as an additive combination of the true grid cell temperature anomaly values at the
 1000 observed grid cells, denoted \mathbf{g}_{obs} , and an observational error term $\boldsymbol{\varepsilon}$:

1001

$$\mathbf{y} = \mathbf{g}_{obs} + \boldsymbol{\varepsilon} \quad (\text{A1})$$

1002

1003 The temperature anomaly field is decomposed into a regression model for the field mean,
 1004 described in terms of a matrix of basis functions \mathbf{H} with coefficients $\boldsymbol{\beta}$, and a spatially correlated
 1005 field \mathbf{f} . The observations are then modelled by this decomposition, notating the basis function
 1006 and the spatial field values at the observed grid cells as \mathbf{H}_{obs} and \mathbf{f}_{obs} :

1007

$$\mathbf{y} = \mathbf{f}_{obs} + \mathbf{H}_{obs}^T \boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (\text{A2})$$

1008

1009 Similarly, we define \mathbf{g}_* as the values true temperature anomaly values at a set of prediction grid
 1010 cells, notating basis functions and the spatial random field values at the prediction grid cells as
 1011 \mathbf{H}_* and \mathbf{f}_* , so that $\mathbf{g}_* = \mathbf{f}_* + \mathbf{H}_*^T \boldsymbol{\beta}$. In this analysis, \mathbf{H} is set as a vector of ones so that the
 1012 regression model acts as an estimate of a constant field mean for the analyzed month.

1013 The spatial field \mathbf{f} is defined in terms of its covariance structure. This covariance structure is
 1014 parameterized as a function of distance between locations, as is common in Gaussian process or
 1015 kriging analyses. The covariance $k(s_m, s_n)$ in spatial field values between locations s_m and s_n is
 1016 defined as:

$$k(s_m, s_n) = \text{cov}(f(s_m), f(s_n)) \quad (\text{A3})$$

1018

1019 which defines the elements of a covariance matrix \mathbf{K} , with elements $[\mathbf{K}]_{mn} = k(s_m, s_n)$. In this
 1020 analysis, a Matérn covariance function is used to model the covariances $k(s_m, s_n)$. This
 1021 covariance function is parameterized by a smoothing hyperparameter ν , a range hyperparameter
 1022 r that controls the rate at which covariance decays with distance between locations, and an
 1023 amplitude hyperparameter σ . We use a stationary covariance function, with fixed values of the
 1024 model hyperparameters fitted independently for the land air temperature and sea-surface
 1025 temperature analyses. Covariances are evaluated as a function of Euclidian distance, rather than
 1026 great circle distance, to retain the flexibility of Matérn covariance functions for data on the
 1027 surface of a spherical Earth, avoiding restrictions to the range of smoothing hyperparameter
 1028 values ν for which Matérn covariances are valid (i.e. to produce positive-definite covariance
 1029 matrices) when using great circle distances (Gneiting, 2013). For separation distances with
 1030 sufficiently strong covariance to be physically important, the Euclidian distance is close to the
 1031 great circle distance.

1032 Values of the field at observed grid cells, \mathbf{f}_{obs} , are modelled as realizations from
 1033 $\mathbf{f}_{obs} \sim N(\mathbf{0}, \mathbf{K}_{obs})$ while those at predictions locations, \mathbf{f}_* , are modelled as $\mathbf{f}_* \sim N(\mathbf{0}, \mathbf{K}_*)$. Cross
 1034 covariances between observed grid cells and prediction grid cells (i.e. the full output grid) are
 1035 defined as \mathbf{K}_{cross} . We define \mathbf{K}_y as the sum of the covariance \mathbf{K}_{obs} and the observational error
 1036 covariance \mathbf{R} :

1037

$$\mathbf{K}_y = \mathbf{K}_{obs} + \mathbf{R} \quad (\text{A4})$$

1038

1039 The observational error covariance matrices are constructed from the error model terms of the
 1040 non-infilled datasets. When the analysis method is applied to an ensemble member of the land air
 1041 temperature ensemble (i.e. the observation vector \mathbf{y} contains the grid cell values for an individual
 1042 land ensemble member for one month), the observational error covariance \mathbf{R} contains the
 1043 additional uncorrelated within-grid-cell measurement and sampling error variances on the
 1044 leading diagonal with zeros elsewhere. When applied to a sea-surface temperature ensemble
 1045 member (i.e. \mathbf{y} contains the grid cell values for an individual HadSST4 ensemble member), \mathbf{R} is
 1046 constructed from the HadSST4 per-platform uncertainties for the partially correlated error
 1047 component, provided as full error covariances in HadSST4, with additional uncertainty from
 1048 uncorrelated measurement and sampling error variances added onto the leading diagonal.

1049 Estimation proceeds following Rasmussen & Williams (2006). The expected value of the
1050 anomaly field \mathbf{g}_* given the observations \mathbf{y} is defined as $\boldsymbol{\mu}_{\mathbf{g}_*|\mathbf{y}}$ where:

1051

$$\boldsymbol{\mu}_{\mathbf{g}_*|\mathbf{y}} = \mathbf{K}_{cross}^T \mathbf{K}_y^{-1} \mathbf{y} + \mathbf{F}^T \boldsymbol{\mu}_{\boldsymbol{\beta}|\mathbf{y}} \quad (\text{A5})$$

1052

1053

1054 and:

1055

$$\mathbf{F}_* = \mathbf{H}_* - \mathbf{H}_{obs} \mathbf{K}_y^{-1} \mathbf{K}_{cross} \quad (\text{A6})$$

1056

1057 Here the terms involving the estimation of regression coefficients $\boldsymbol{\beta}$ (of which we need no prior
1058 knowledge) are:

1059

$$\boldsymbol{\mu}_{\boldsymbol{\beta}|\mathbf{y}} = \boldsymbol{\Sigma}_{\boldsymbol{\beta}|\mathbf{y}} \mathbf{H}_{obs} \mathbf{K}_y^{-1} \mathbf{y} \quad (\text{A7})$$

1060

$$\boldsymbol{\Sigma}_{\boldsymbol{\beta}|\mathbf{y}} = (\mathbf{H}_{obs} \mathbf{K}_y^{-1} \mathbf{H}_{obs}^T)^{-1} \quad (\text{A8})$$

1061

1062 The posterior covariance $\boldsymbol{\Sigma}_{\mathbf{g}_*|\mathbf{y}}$ for the Gaussian process prediction is given by:

1063

$$\boldsymbol{\Sigma}_{\mathbf{g}_*|\mathbf{y}} = \mathbf{K}_* - \mathbf{K}_{cross}^T \mathbf{K}_y \mathbf{K}_{cross} + \mathbf{F}^T \boldsymbol{\Sigma}_{\boldsymbol{\beta}|\mathbf{y}} \mathbf{F} \quad (\text{A9})$$

1064

1065 Together, $\boldsymbol{\mu}_{\mathbf{g}_*|\mathbf{y}}$ and $\boldsymbol{\Sigma}_{\mathbf{g}_*|\mathbf{y}}$ define the full posterior distribution of the Gaussian process estimate
1066 of the gridded temperature anomaly field \mathbf{g}_* for all output grid cells, given observations \mathbf{y} .

1067 *A.2 Kernel hyperparameter estimation*

1068 The estimation of the amplitude (σ) and decorrelation range (r) parameters of our spatial model
1069 is based on application of the maximum marginal likelihood method that is described in
1070 Rasmussen & Williams (2006). Here, the kernel hyperparameters $\boldsymbol{\theta} = (\sigma, r)$ are fit through

1071 numerical optimization to find the parameters that maximize the marginal log likelihood
 1072 function, rearranged here as:

1073

$$\log p(\mathbf{y}|\boldsymbol{\theta}) = -\frac{1}{2}\mathbf{y}^T\mathbf{K}_y^{-1}\mathbf{y} + \frac{1}{2}\boldsymbol{\mu}_{\beta|\mathbf{y}}^T\boldsymbol{\Sigma}_{\beta|\mathbf{y}}^{-1}\boldsymbol{\mu}_{\beta|\mathbf{y}} - \frac{1}{2}\log|\mathbf{K}_y| + \frac{1}{2}\log|\boldsymbol{\Sigma}_{\beta|\mathbf{y}}| - \frac{N-J}{2}\log(2\pi) \quad (\text{A1 } 0)$$

1074

1075 Here, N is the number of observed grid cells in \mathbf{y} and J is the number of covariates included in
 1076 the regression portion of the analysis model. We include a single covariate for the analysis field
 1077 mean, hence $J = 1$ in our application.

1078 The hyperparameters are fit to monthly ‘best estimate’ gridded temperature anomaly fields
 1079 separately for land air temperatures and sea-surface temperatures. Observational uncertainties are
 1080 derived from the HadCRUT5 land ensemble uncertainty model (described in Morice et al., 2012)
 1081 and HadSST4 uncertainty model (Kennedy et al., 2019), as described below.

1082 As we fit hyperparameters to best estimates of the non-filled grids, we include an additional
 1083 uncertainty component in the observational error covariance to represent the observational bias
 1084 uncertainty that is encoded into the land ensemble and the HadSST4 ensemble. Hence, when
 1085 fitting hyperparameters, an extended observational error covariance \mathbf{R}' is substituted for \mathbf{R} where
 1086 $\mathbf{R}' = \mathbf{R} + \boldsymbol{\Sigma}_{ensemble}$ and $\boldsymbol{\Sigma}_{ensemble}$ is an error covariance matrix that is empirically derived from
 1087 the ensemble. The ensemble-derived error covariance matrices are only used when fitting
 1088 hyperparameters for the best estimate fields. They are not included in the observational error
 1089 covariance term when fitting the analysis fields for individual ensemble members in Appendix
 1090 A.3.

1091 For land hyperparameter estimation, the monthly observation vector \mathbf{y} is constructed from a
 1092 CRUTEM5 best estimate field. The observational error covariance \mathbf{R} is constructed from the
 1093 uncorrelated measurement and sampling uncertainty grids, from the Brohan et al. (2006) error
 1094 model, while $\boldsymbol{\Sigma}_{ensemble}$ is computed from the HadCRUT5 land ensemble. For marine
 1095 hyperparameter estimation, the observation vector \mathbf{y} is constructed from a HadSST4 ensemble
 1096 median field. The observational error covariance matrices \mathbf{R} are constructed by combining
 1097 HadSST4 uncorrelated measurement and sampling uncertainties with the HadSST4 ‘micro bias’
 1098 error covariance matrices and $\boldsymbol{\Sigma}_{ensemble}$ is computed from the HadSST4 ensemble.

1099 Hyperparameter estimates are computed for each of the 360 monthly fields in the 1961 to 1990
 1100 climatology period, during which the observational sampling is near global in extent. The
 1101 hyperparameters used in the analysis are taken as the average of the hyperparameters fitted in the
 1102 360 monthly optimizations, with scale parameters rounded to the nearest 0.05 °C and range
 1103 parameters rounded to the nearest 50 km. The resulting amplitude parameter σ and range
 1104 parameter r for the land air temperature analysis are $\sigma = 1.2^\circ\text{C}$ and $r = 1300$ km. For the sea
 1105 surface temperature analysis, the fitted parameters are $\sigma = 0.6^\circ\text{C}$ and $r = 1300$ km. The
 1106 smoothing parameter was fixed at $\nu = 1.5$. This model represents typical land and marine

1107 temperature anomaly variability. The model does not include regional and seasonal variations in
 1108 these parameters, nonetheless where there is a sufficient observational constraint the method can
 1109 reproduce appropriate regional and seasonal variability in the analysis anomaly fields. Additional
 1110 information on the monthly hyperparameter fits can be found in the Supporting Information.

1111 *A.3 Ensemble analysis*

1112 The HadCRUT5 ensemble land and marine analyses are constructed by applying Gaussian
 1113 process regression to each ensemble member of the non-infilled land and marine data sets.
 1114 Uncertainty is further explored by encoding analysis uncertainty into the ensemble, sampling
 1115 from the Gaussian process posterior distribution through a process called conditional simulation
 1116 (Chilès & Delfiner, 2012).

1117 We denote a vector of observed grid cell temperature anomalies for a non-infilled ensemble
 1118 member as \mathbf{y}_d , with the subscript d indexing the ensemble member. We then apply the Gaussian
 1119 process analysis method to compute the expected value of the temperature anomaly field $\boldsymbol{\mu}_{\mathbf{g}_*|y_d}$
 1120 for the ensemble member, substituting \mathbf{y}_d and $\boldsymbol{\mu}_{\mathbf{g}_*|y_d}$ into Equation A5. We then proceed to
 1121 sample the analysis uncertainty through conditional simulation, as described below.

1122 For each ensemble member, we draw a random sample from the joint prior distribution of the
 1123 anomaly field at observed and prediction locations, setting the regression coefficient for each
 1124 sample to an arbitrary value of $\boldsymbol{\beta}' = \mathbf{0}$. This sampling distribution is defined as:

1125

$$\begin{bmatrix} \mathbf{g}'_{obs} \\ \mathbf{g}'_* \end{bmatrix} \sim N \left(\begin{bmatrix} \mathbf{H}_{obs}^T \mathbf{0} \\ \mathbf{H}_*^T \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{obs} & \mathbf{K}_{cross}^T \\ \mathbf{K}_{cross} & \mathbf{K}_* \end{bmatrix} \right) \quad (\text{A11})$$

1126

1127 This provides samples of the anomaly field, according to the Gaussian process model on the full
 1128 output grid, drawn as $\mathbf{g}'_* = \mathbf{f}'_* + \mathbf{H}_*^T \mathbf{0}$, and at the observed locations $\mathbf{g}'_{obs} = \mathbf{f}'_{obs} + \mathbf{H}_{obs}^T \mathbf{0}$, with
 1129 the correct covariance structure between observed and output grid locations.

1130 We then generate pseudo-observations \mathbf{y}' of the simulated temperature field by sampling from
 1131 the observational error model $\boldsymbol{\varepsilon}' \sim N(\mathbf{0}, \mathbf{R})$. The simulated observation is then defined as:

1132

$$\mathbf{y}' = \mathbf{f}'_{obs} + \mathbf{H}_{obs}^T \mathbf{0} + \boldsymbol{\varepsilon}' \quad (\text{A12})$$

1133

1134 Simulations of reconstruction error are based on application of the Gaussian process estimation
 1135 to the simulated anomaly fields and simulated (pseudo) observations. The difference between the
 1136 simulated field sample \mathbf{g}'_* and the estimate based on the simulated pseudo observations $\boldsymbol{\mu}_{\mathbf{g}'_*|y'}$ is
 1137 a sample of the reconstruction error according to the Gaussian process model. This difference,

1138 $\mathbf{e}' = \boldsymbol{\mu}_{\mathbf{g}'_*|\mathbf{y}'} - \mathbf{g}'_*$, is a sample from the posterior distribution of the Gaussian process regression,
 1139 i.e. $\mathbf{e}' \sim N(\boldsymbol{\mu}_{\mathbf{g}'_*|\mathbf{y}'}, \boldsymbol{\Sigma}_{\mathbf{g}'_*|\mathbf{y}'})$.

1140 For an ensemble member indexed by d with observation vector \mathbf{y}_d , the analysis values \mathbf{g}_{*d} are
 1141 computed as the sum of the Gaussian process estimate $\boldsymbol{\mu}_{\mathbf{g}_{*d}|\mathbf{y}_d}$, based on the real observations
 1142 \mathbf{y}_d , and a simulated reconstruction error sample \mathbf{e}'_d :

1143

$$\mathbf{g}_{*d} = \boldsymbol{\mu}_{\mathbf{g}_{*d}|\mathbf{y}_d} + \mathbf{e}'_d \quad (\text{A13})$$

1144

1145 The resulting ensemble encodes both the bias terms in the underlying observational ensemble
 1146 and the reconstruction error for the Gaussian process.

1147 The applied Gaussian process estimation is purely spatial and so does not provide information on
 1148 temporally-correlated reconstruction error. To mitigate this, we modify the above sampling
 1149 method to encode temporal correlation into the conditional simulation process. The simulated
 1150 spatial fields \mathbf{g}'_* and \mathbf{g}'_{obs} are sampled such that they are fully correlated throughout a year, i.e.
 1151 the same spatial field is used for each sample within a year. This provides a conservative upper
 1152 bound on uncertainty in annual averages derived from the ensemble.

1153 Known temporal correlations in observational measurement and sampling errors, which are not
 1154 represented in the non-infilled land and marine ensembles, are similarly encoded into the
 1155 observational error samples $\boldsymbol{\varepsilon}'_d$ when generating pseudo-observations. This strategy is applied for
 1156 the residual SST micro biases that are represented in the HadSST4 observational error
 1157 covariance matrices. These are encoded using the same random draw for all months in a year
 1158 when sampling. This allows uncertainty in annual averages to be computed under a conservative
 1159 assumption of full temporal correlation of SST micro biases within a year, as defined by the
 1160 HadSST4 uncertainty model (Kennedy et al., 2019). Other measurement and sampling
 1161 uncertainties, associated with temporally uncorrelated errors, are sampled independently for each
 1162 month. No additional temporal correlation is encoded into the ensemble for land air temperatures
 1163 as there is no temporal correlation in the measurement and sampling error terms for CRUTEM5
 1164 (although the analyzed land ensemble does already sample time correlated observational errors
 1165 from residual station biases, which are distinct from the measurement and sampling uncertainty
 1166 terms discussed here).

1167 Although knowledge of temporal correlation in errors is not used to improve the estimated
 1168 anomaly fields, the result of the sampling process is to enable an upper bound on uncertainty in
 1169 annual averages to be obtained directly from the ensemble.

1170 *A.4 Observational constraint mask*

1171 Despite the application of spatial reconstruction, there are regions of the world in which the
 1172 available observational coverage, particularly in the early part of the record, is such that a
 1173 reliable reconstruction is not possible. In regions where local observations are not available, the

1174 analysis ensemble mean reverts towards the regression model estimate of the mean temperature
 1175 anomaly, inferred from observed regions, while the ensemble spread tends towards that
 1176 described by the Gaussian process prior distribution.

1177 Consequently, regions where the constraint from local observations is poor are removed from the
 1178 analysis. The reconstruction in these regions is highly sensitive to the prior covariance model and
 1179 the estimated regression term $\mathbf{H}_*^T \boldsymbol{\mu}_{\beta|y}$, for which the coefficient estimate may be biased towards
 1180 observed regions. This has been found to be the case in test analyses of climate model
 1181 simulations in which global average temperature estimates have been found to be biased towards
 1182 northern hemisphere temperatures during periods with sparse southern hemisphere coverage.

1183 The criteria used to mask regions, defined in terms of a threshold α , is based on the ratio of
 1184 posterior and prior variance of the local Gaussian process estimate, omitting the global
 1185 regression term which has an improper prior, with regions of the analysis masked where the
 1186 following inequality is satisfied:

1187

$$\mathbf{1} - \frac{\text{diag}(\mathbf{K}_* - \mathbf{K}_{cross}^T \mathbf{K}_y \mathbf{K}_{cross})}{\text{diag}(\mathbf{K}_*)} < \alpha \quad (\text{A14})$$

1188

1189 The left-hand side of Equation A14 is bounded between zero and one and we use a threshold of
 1190 $\alpha = 0.25$ to provide a balance between retaining regions with useful information content and
 1191 masking those regions that have a weak observational constraint. Global and hemispheric
 1192 average temperature series for varying α are provided in the Supporting Information and indicate
 1193 that these diagnostics are insensitive to the choice of α values in the range 0.1 to 0.5.

1194

1195 **References**

- 1196 Allen, M. R., O. P. Dube, W. Solecki, F. Arag3n-Durand, W. Cramer, S. Humphreys, M.
 1197 Kainuma, J. Kala, N. Mahowald, Y. Mulugetta, R. Perez, M. Wairiu, and K. Zickfeld
 1198 (2018), Framing and Context. In: *Global Warming of 1.5°C. An IPCC Special Report on*
 1199 *the impacts of global warming of 1.5°C above pre-industrial levels and related global*
 1200 *greenhouse gas emission pathways, in the context of strengthening the global response to*
 1201 *the threat of climate change, sustainable development, and efforts to eradicate poverty*
 1202 *[Masson-Delmotte, V., P. Zhai, H.-O. P3rtner, D. Roberts, J. Skea, P.R. Shukla, A.*
 1203 *Pirani, W. Moufouma-Okia, C. P3an, R. Pidcock, S. Connors, J.B.R. Matthews, Y. Chen,*
 1204 *X. Zhou, M.I. Gomis, E. Lonnoy, T. Maycock, M. Tignor, and T. Waterfield (eds.)].*
- 1205 Atkinson, C. P., N. A. Rayner, J. J. Kennedy and S. A. Good (2014), An integrated database of
 1206 ocean temperature and salinity observations, *Journal of Geophysical Research: Oceans*,
 1207 119, 7139– 7163, doi:10.1002/2014JC010053

- 1208 Benestad, R. E., H. B. Erlandsen, A. Mezghani and K. M. Parding (2019), Geographical
 1209 distribution of thermometers gives the appearance of lower historical global warming,
 1210 *Geophysical Research Letters*, 46. <https://doi.org/10.1029/2019GL083474>
- 1211 Berry, D. I., E. C. Kent, and P. K. Taylor (2004), An Analytical Model of Heating Errors in
 1212 Marine Air Temperatures from Ships, *J. Atmos. Oceanic Technol.*, 21, 1198–1215,
 1213 doi:10.1175/1520-0426(2004)021<1198:AAMOHE>2.0.CO;2
- 1214 Berry, D. I. and Kent, E. C. (2017), Assessing the health of the in situ global surface marine
 1215 climate observing system. *Int. J. Climatol.*, 37: 2248-2259. doi:10.1002/joc.4914
- 1216 Blunden, J. and D. S. Arndt, Eds. (2019), State of the Climate in 2018. *Bull. Amer. Meteor. Soc.*,
 1217 100 (9), Si–S305, doi:10.1175/2019BAMSSStateoftheClimate.1.
- 1218 Brohan, P., J. J. Kennedy, I. Harris, S. F. B. Tett and P. D. Jones (2006), Uncertainty estimates in
 1219 regional and global observed temperature changes: a new dataset from 1850, *J. Geophys.*
 1220 *Res*, 111, D12106, doi:10.1029/2005JD006548.
- 1221 Carella, G., E. C. Kent, and D. I. Berry (2017), A probabilistic approach to ship voyage
 1222 reconstruction in ICOADS, *Int. J. Climatol.*, 37, 2233-2247, doi:10.1002/joc.4492
- 1223 Carella, G., J. J. Kennedy, D. I. Berry, S. Hirahara, C. J. Merchant, S. Morak-Bozzo and E. C.
 1224 Kent (2018), Estimating sea surface temperature measurement methods using
 1225 characteristic differences in the diurnal cycle. *Geophysical Research Letters*, 45, 363–
 1226 371. <https://doi.org/10.1002/2017GL076475>
- 1227 Chan, D. and P. Huybers (2019), Systematic Differences in Bucket Sea Surface Temperature
 1228 Measurements among Nations Identified Using a Linear-Mixed-Effect Method, *J.*
 1229 *Climate*, 32, 2569–2589, doi:10.1175/JCLI-D-18-0562.1
- 1230 Chilès, J-P and P. Delfiner (2012). *Geostatistics: Modeling Spatial Uncertainty*. Wiley Series In
 1231 *Probability and Statistics*. 10.1002/9781118136188.
- 1232 Cornes, R. C., E. Kent, D. Berry and J. J. Kennedy (2020), CLASSnmat: A global night marine
 1233 air temperature data set, 1880–2019. *Geosci Data J.*, 7, 170–184.
 1234 <https://doi.org/10.1002/gdj3.100>
- 1235 Cowtan, K. and R. G. Way (2014), Coverage bias in the HadCRUT4 temperature series and its
 1236 impact on recent temperature trends, *Q.J.R. Meteorol. Soc.*, 140, 1935–1944.
 1237 doi:10.1002/qj.2297
- 1238 Cowtan, K., Z. Hausfather, E. Hawkins, P. Jacobs, M. E. Mann, S. K. Miller, B. A. Steinman, M.
 1239 B. Stolpe and R. G. Way (2015), Robust comparison of climate models with observations
 1240 using blended land air and ocean sea surface temperatures, *Geophys. Res. Lett.*, 42,
 1241 6526– 6534, doi:10.1002/2015GL064888.
- 1242 Cowtan, K., R. Rohde, Z. Hausfather (2018), Evaluating biases in sea surface temperature
 1243 records using coastal weather stations, *Q J R Meteorol Soc.*, 144, 670–681,
 1244 doi:10.1002/qj.3235
- 1245 Donlon, C. J., M. Martin, J. Stark, J. Roberts-Jones, E. Fiedler and W. Wimmer (2012), The
 1246 operational sea surface temperature and sea ice analysis (OSTIA) system, *Remote*
 1247 *Sensing of Environment*, 116, 140-158, doi:10.1016/j.rse.2010.10.017

- 1248 Folland, C. K. and D. E. Parker (1995), Correction of instrumental biases in historical sea surface
1249 temperature data, *Quarterly Journal of the Royal Meteorological Society*, 121, 522, 319–
1250 367, doi:10.1002/qj.49712152206.
- 1251 Folland, C. K., N. A. Rayner, S. J. Brown, T. M. Smith, S. S. P. Shen, D. E. Parker, I. Macadam,
1252 P. D. Jones, R. N. Jones, N. Nichols and D. M. H. Sexton (2001), Global temperature
1253 change and its uncertainties since 1861, *Geophysical Research Letters*, 28(13), 2621–
1254 2624
- 1255 Freeman, E., S. D. Woodruff, S. J. Worley, S. J. Lubker, E.C. Kent, W. E. Angel, D. I. Berry, P.
1256 Brohan, R. Eastman, L. Gates, W. Gloeden, Z. Ji, J. Lawrimore, N.A. Rayner, G.
1257 Rosenhagen and S. R. Smith (2017), ICOADS Release 3.0: a major update to the
1258 historical marine climate record. *Int. J. Climatol.*, 37, 2211–2232, doi:10.1002/joc.4775.
- 1259 Gelaro, R., W. McCarty, M.J. Suárez, R. Todling, A. Molod, L. Takacs, C.A. Randles, A.
1260 Darnenov, M. G. Bosilovich, R. Reichle, K. Wargan, L. Coy, R. Cullather, C. Draper, S.
1261 Akella, V. Buchard, A. Conaty, A.M. da Silva, W. Gu, G. Kim, R. Koster, R. Lucchesi,
1262 D. Merkova, J. E. Nielsen, G. Partyka, S. Pawson, W. Putman, M. Rienecker, S. D.
1263 Schubert, M. Sienkiewicz, and B. Zhao, 2017: The Modern-Era Retrospective Analysis
1264 for Research and Applications, Version 2 (MERRA-2). *J. Climate*, 30, 5419–5454,
1265 <https://doi.org/10.1175/JCLI-D-16-0758.1>
- 1266 Gneiting, T. (2013). Strictly and non-strictly positive definite functions on spheres, *Bernoulli*, 19,
1267 1327–1349, doi:10.3150/12-BEJSP06
- 1268 Hansen, J., R. Ruedy, M. Sato, and K. Lo (2010), Global surface temperature change, *Rev.*
1269 *Geophys.*, 48, RG4004, doi:10.1029/2010RG000345.
- 1270 Hartmann, D. L., A. M. G. Klein Tank, M. Rusticucci, L.V. Alexander, S. Brönnimann, Y.
1271 Charabi, F. J. Dentener, E. J. Dlugokencky, D. R. Easterling, A. Kaplan, B. J. Soden, P.
1272 W. Thorne, M. Wild and P. M. Zhai (2013), Observations: Atmosphere and Surface. In:
1273 *Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to*
1274 *the Fifth Assessment Report of the Intergovernmental Panel on Climate Change* [Stocker,
1275 T.F., D. Qin, G.-K. Plattner, M. Tignor, S.K. Allen, J. Boschung, A. Nauels, Y. Xia, V.
1276 Bex and P.M. Midgley (eds.)]. Cambridge University Press, Cambridge, United Kingdom
1277 and New York, NY, USA.
- 1278 Hausfather, Z., K. Cowtan, D.C. Clarke, P. Jacobs, M. Richardson, and R. Rohde (2017),
1279 Assessing recent warming using instrumentally homogeneous sea surface temperature
1280 records, *Science advances*, 3, 1, p.e1601207, doi:10.1126/sciadv.1601207.
- 1281 Hawkins, E. et al., (2017), Estimating changes in global temperature since the pre-industrial
1282 period. *Bulletin of the American Meteorological Society*, BAMS–D–16–0007.1,
1283 doi:10.1175/bams-d-16-0007.1.
- 1284 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., Nicolas, J.,
1285 Peubey, C., Radu, R., Schepers, D., Simmons, A., Soci, C., Abdalla, S., Abellan, X.,
1286 Balsamo, G., Bechtold, P., Biavati, G., Bidlot, J., Bonavita, M., De Chiara, G., Dahlgren,
1287 P., Dee, D., Diamantakis, M., Dragani, R., Flemming, J., Forbes, R., Fuentes, M., Geer,
1288 A., Haimberger, L., Healy, S., Hogan, R.J., Hólm, E., Janisková, M., Keeley, S.,
1289 Laloyaux, P., Lopez, P., Lupu, C., Radnoti, G., de Rosnay, P., Rozum, I., Vamborg, F.,
1290 Villaume, S. and Thépaut, J.-N. (2020), The ERA5 Global Reanalysis. *Q J R Meteorol*

- 1291 Soc. Accepted Author Manuscript. doi:10.1002/qj.3803 Hirahara, S., M. Ishii, and Y.
 1292 Fukuda (2014), Centennial-scale sea surface temperature analysis and its uncertainty, *J.*
 1293 *Climate*, 27, 57–75, doi:10.1175/JCLI-D-12-00837.1,
- 1294 Huang, B., P. W. Thorne, T. M. Smith, W. Liu, J. H. Lawrimore, V. F. Banzon, H. Zhang, T. C.
 1295 Peterson, M. J. Menne (2016), Further exploring and quantifying uncertainties for
 1296 Extended Reconstructed Sea Surface Temperature (ERSST) version 4 (v4). *J. Climate*,
 1297 doi:10.1175/JCLI-D-15-0430.1.
- 1298 Huang, B., P. W. Thorne, V. F. Banzon, T. Boyer, G. Chepurin, J. H. Lawrimore, M. J. Menne,
 1299 T. M. Smith, R. S. Vose, and H. Zhang (2017), Extended Reconstructed Sea Surface
 1300 Temperature, Version 5 (ERSSTv5): Upgrades, Validations, and Intercomparisons, *J.*
 1301 *Climate*, 30, 8179–8205, doi:10.1175/JCLI-D-16-0836.1.
- 1302 Huang, B., M. J. Menne, T. Boyer, E. Freeman, B. E. Gleason, J. H. Lawrimore, C. Liu, J. J.
 1303 Rennie, C. J. Schreck, F. Sun, R. Vose, C. N. Williams, X. Yin, and H. Zhang (2019),
 1304 Uncertainty estimates for sea surface temperature and land surface air temperature in
 1305 NOAA GlobalTemp version 5, *J. Climate*, 0, doi:10.1175/JCLI-D-19-0395.1
- 1306 Ilyas, M., C. M. Brierley and S. Guillas (2017), Uncertainty in regional temperatures inferred
 1307 from sparse global observations: Application to a probabilistic classification of El Niño,
 1308 *Geophys. Res. Lett.*, 44, 9068– 9074, doi:10.1002/2017GL074596.
- 1309 Jones, P. D., D. H. Lister, T. J. Osborn, C. Harpham, M. Salmon, and C. P. Morice (2012),
 1310 Hemispheric and large-scale land surface air temperature variations: An extensive
 1311 revision and an update to 2010, *J. Geophys. Res.*, 117, D05127,
 1312 doi:10.1029/2011JD017139
- 1313 Jones, P. D., T. J. Osborn, and K. R. Briffa (1997), Estimating Sampling Errors in Large-Scale
 1314 Temperature Averages. *J. Climate*, 10, 2548–2568, [https://doi.org/10.1175/1520-0442\(1997\)010<2548:ESEILS>2.0.CO;2](https://doi.org/10.1175/1520-0442(1997)010<2548:ESEILS>2.0.CO;2)
- 1316 Junod, R. A. and J. R. Christy (2020), A new compilation of globally gridded night-time marine
 1317 air temperatures: The UAHNMATv1 dataset. *Int J Climatol*, 40, 2609– 2623.
 1318 <https://doi.org/10.1002/joc.6354>
- 1319 Kadow, C., Hall, D. M. & Ulbrich, U. (2020) Artificial intelligence reconstructs missing climate
 1320 information. *Nat. Geosci.*, <https://doi.org/10.1038/s41561-020-0582-5>
- 1321 Kalnay, E., M. Kanamitsu, R. Kistler, W. Collins, D. Deaven, L. Gandin, M. Iredell, S. Saha, G.
 1322 White, J. Woollen, Y. Zhu, M. Chelliah, W. Ebisuzaki, W. Higgins, J. Janowiak, K.C.
 1323 Mo, C. Ropelewski, J. Wang, A. Leetmaa, R. Reynolds, R. Jenne, and D. Joseph (1996),
 1324 The NCEP/NCAR 40-Year Reanalysis Project. *Bull. Amer. Meteor. Soc.*, 77, 437–472,
 1325 [https://doi.org/10.1175/1520-0477\(1996\)077<0437:TNYRP>2.0.CO;2](https://doi.org/10.1175/1520-0477(1996)077<0437:TNYRP>2.0.CO;2)
- 1326 Karl, T. R., A. Arguez, B. Huang, J. H. Lawrimore, J. R. McMahon, M. J. Menne, T. C.
 1327 Peterson, R. S. Vose, and H-M. Zhang (2015), Possible artifacts of data biases in the
 1328 recent global surface warming hiatus, *Science*, 348(6242): 1469–1472,
 1329 doi:10.1126/science.aaa5632.
- 1330 Kaplan, A., Y. Kushnir, M. A. Cane, and M. B. Blumenthal (1997), Reduced space optimal
 1331 analysis for historical data sets: 136 years of Atlantic sea surface temperatures, *J.*
 1332 *Geophys. Res.*, 102(C13), 27835– 27860, doi:10.1029/97JC01734.

- 1333 Kennedy J. J., N. A. Rayner, R. O. Smith, M. Saunby and D. E. Parker (2011a), Reassessing
 1334 biases and other uncertainties in sea-surface temperature observations since 1850 part 1:
 1335 measurement and sampling errors, *J. Geophys. Res.*, 116, D14103,
 1336 doi:10.1029/2010JD015218
- 1337 Kennedy J. J., N. A. Rayner, R. O. Smith, M. Saunby and D. E. Parker (2011b), Reassessing
 1338 biases and other uncertainties in sea-surface temperature observations since 1850 part 2:
 1339 biases and homogenisation, *J. Geophys. Res.*, 116, D14104, doi:10.1029/2010JD015220
- 1340 Kennedy, J. J., N. A. Rayner, C. P. Atkinson, and R. E. Killick (2019), An ensemble data set of
 1341 sea-surface temperature change from 1850: the Met Office Hadley Centre
 1342 HadSST.4.0.0.0 data set, *Journal of Geophysical Research: Atmospheres*, 124,
 1343 doi:10.1029/2018JD029867
- 1344 Kent, E. C., N. A. Rayner, D. I. Berry, M. Saunby, B. I. Moat, J. J. Kennedy and D. E. Parker
 1345 (2013), Global analysis of night marine air temperature and its uncertainty since 1880:
 1346 The HadNMAT2 data set, *J. Geophys. Res. Atmos.*, 118, 1281–1298,
 1347 doi:10.1002/jgrd.50152.
- 1348 Kent, E. C., J. J. Kennedy, T. M. Smith, S. Hirahara, B. Huang, A. Kaplan, D. E. Parker, C. P.
 1349 Atkinson, D. I. Berry, G. Carella, Y. Fukuda, M. Ishii, P. D. Jones, F. Lindgren, C. J.
 1350 Merchant, S. Morak-Bozzo, N. A. Rayner, V. Venema, S. Yasui, and H. Zhang (2017), A
 1351 Call for New Approaches to Quantifying Biases in Observations of Sea Surface
 1352 Temperature, *Bull. Amer. Meteor. Soc.*, 98, 1601–1616, doi:10.1175/BAMS-D-15-
 1353 00251.1.
- 1354 Klein Tank, A. M. G., J. B. Wijngaard, G.P.Können, R. Böhm, G. Demarée, A. Gocheva, M.
 1355 Mileta, S. Pashiardis, L. Hejkrlik, C. Kern-Hansen, R. Heino, P. Bessemoulin, G. Müller-
 1356 Westermeier, M. Tzanakou, S. Szalai, T. Pálsdóttir, D. Fitzgerald, S. Rubin, M. Capaldo,
 1357 M., Maugeri, A. Leitass, A. Bukantis, R. Aberfeld, A. F. V. van Engelen, E. Forland, M.
 1358 Mietus, F. Coelho, C. Mares, V. Razuvaev, E. Nieplova, T. Cegnar, J. Antonio López, B.
 1359 Dahlström, A. Moberg, W. Kirchhofer, A. Ceylan, O. Pachaliuk, L. V. Alexander and P.
 1360 Petrovic (2002), Daily dataset of 20th-century surface air temperature and precipitation
 1361 series for the European Climate Assessment. *Int. J. Climatol.*, 22: 1441-1453.
 1362 doi:10.1002/joc.773
- 1363 Kobayashi, S., et al. (2015), The JRA-55 Reanalysis: General Specifications and Basic
 1364 Characteristics, *Journal of the Meteorological Society of Japan. Ser. II*, 93(1), 5-48,
 1365 doi:10.2151/jmsj.2015-001.
- 1366 Lenssen, N., G. Schmidt, J. Hansen, M. Menne, A. Persin, R. Ruedy, and D. Zyss (2019),
 1367 Improvements in the GISTEMP uncertainty model, *J. Geophys. Res. Atmos.*, 124, 12,
 1368 6307-6326, doi:10.1029/2018JD029522.
- 1369 Menne, M. J., C. N. Williams, B. E. Gleason, J. J. Rennie, and J. H. Lawrimore (2018), The
 1370 Global Historical Climatology Network Monthly Temperature Dataset, Version 4, *J.*
 1371 *Climate*, 31, 9835–9854, doi:10.1175/JCLI-D-18-0094.1.
- 1372 Osborn, T. J., P. D. Jones, D. H. Lister, C. P. Morice, I. R. Simpson and I. C. Harris (2020), Land
 1373 surface air temperature variations across the globe updated to 2019: the CRUTEM5
 1374 dataset, Submitted to *J. Geophys. Res.*

- 1375 Morice, C. P., J. J. Kennedy, N. A. Rayner, and P. D. Jones (2012), Quantifying uncertainties in
 1376 global and regional temperature change using an ensemble of observational estimates:
 1377 The HadCRUT4 dataset, *J. Geophys. Res.*, 117, D08101, doi:10.1029/2011JD017187.
- 1378 Parker, D.E. (1994), Effects of changing exposure of thermometers at land stations. *Int. J.*
 1379 *Climatol.*, 14: 1-31. doi:10.1002/joc.3370140102
- 1380 Parker, D. E., Jones, P. D., Folland, C. K., and Bevan, A. (1994), Interdecadal changes of surface
 1381 temperature since the late nineteenth century, *J. Geophys. Res.*, 99 (D7), 14373– 14399,
 1382 doi:10.1029/94JD00548
- 1383 Parker, D.E. (2006), A Demonstration That Large-Scale Warming Is Not Urban. *J. Climate*, 19,
 1384 2882–2895, <https://doi.org/10.1175/JCLI3730.1>
- 1385 Parker, D. E. (2010), Urban heat island effects on estimates of observed climate change, *Wiley*
 1386 *Interdisciplinary Reviews: Climate Change*, 1(1), 123-133
- 1387 Rasmussen, C. E. and C. K. I. Williams (2006), *Gaussian Processes for Machine Learning*, the
 1388 MIT Press, 2006, ISBN 026218253X.
- 1389 Rayner, N. A., P. Brohan, D. E. Parker, C. K. Folland, J. J. Kennedy, M. Vanicek, T. J. Ansell
 1390 and S. F. B. Tett (2006), Improved analyses of changes and uncertainties in sea surface
 1391 temperature measured in situ since the mid-nineteenth century: the HadSST2 data set,
 1392 *Journal of Climate*, 19, 3, 446– 469, doi:10.1175/JCLI3637.1.
- 1393 Rayner, N. A., R. Auchmann, J. Bessembinder, S. Brönnimann, Y. Brugnara, F. Capponi, L.
 1394 Carrea, E. M. A. Dodd, D. Ghent, E. Good, J. J. Kennedy, E. C. Kent, R. E. Killick, P.
 1395 van der Linden, F. Lindgren, K. S. Madsen, C. J. Merchant, H. R. Mitchelson, C. P.
 1396 Morice, P. Nielsen-Englyst, P. F. Ortiz, J. J. Remedios, G. van der Schrier, A. A. Squintu,
 1397 A. Stephens, P. W. Thorne, R. T. Tonboe, T. Trent, K. L. Veal, A. M. Waterfall, K.
 1398 Winfield, J. P. Winn, R. I. Woolway (2020), The EUSTACE project: delivering global
 1399 daily information on surface air temperature, *Bull. Amer. Met. Soc.*,
 1400 <https://doi.org/10.1175/BAMS-D-19-0095.1>.
- 1401 Rennie, J. J., J. H. Lawrimore, B. E. Gleason, P. W. Thorne, C. P. Morice, M. J. Menne, C. N.
 1402 Williams, .N., W. G. de Almeida, J. Christy, M. Flannery, M. Ishihara, K. Kamiguchi, A.
 1403 M. G. Klein-Tank, A. Mhanda, D. H. Lister, V. Razuvaev, M. Renom, M. Rusticucci, J.
 1404 Tandy, S. J. Worley, V. Venema, W. Angel, M. Brunet, B. Dattore, H. Diamond, M. A.
 1405 Lazzara, F. Le Blancq, J. Luterbacher, H. Mächel, J. Revadekar, R. S. Vose, and X. Yin
 1406 (2014), The international surface temperature initiative global land surface databank:
 1407 monthly temperature data release description and methods. *Geosci. Data J.*, 1: 75-102.
 1408 doi:10.1002/gdj3.8
- 1409 Reynolds, R. W. and T. M. Smith (1994), Improved Global Sea Surface Temperature Analyses
 1410 Using Optimum Interpolation. *J. Climate*, 7, 929–948, [https://doi.org/10.1175/1520-0442\(1994\)007<0929:IGSSTA>2.0.CO;2](https://doi.org/10.1175/1520-0442(1994)007<0929:IGSSTA>2.0.CO;2)
- 1411
 1412 Reynolds, R. W., N. A. Rayner, T. M. Smith, D. C. Stokes, and W. Wang (2002), An improved
 1413 in situ and satellite SST analysis for climate. *J. Climate*, 15, 1609-1625.
- 1414 Richardson, M., Cowtan, K., Hawkins, E., & Stolpe, M. B. (2016), Reconciled climate response
 1415 estimates from climate models and the energy budget of Earth. *Nature Climate Change*,
 1416 6(10), 931-935. <https://doi.org/10.1038/nclimate3066>

- 1417 Richardson, M., K. Cowtan, and R.J. Millar (2018), Global temperature definition affects
 1418 achievement of long-term climate goals. *Environmental Research Letters*, 13(5), 054004,
 1419 doi:10.1088/1748-9326/aab305.
- 1420 Rohde, R., R. A. Muller, R. Jacobsen, E. Muller, S. Perlmutter, A. Rosenfeld, J. Wurtele, D.
 1421 Groom and C. Wickham (2013a), A New Estimate of the Average Earth Surface Land
 1422 Temperature Spanning 1753 to 2011, *Geoinfor Geostat: An Overview*, 1:1,
 1423 doi:10.4172/gigs.1000101.
- 1424 Rohde R., R. M. Muller, R. Jacobsen, S. Perlmutter, A. Rosenfeld, J. Wurtele, J. Curry, C.
 1425 Wickham and S. Mosher (2013b), Berkeley Earth Temperature Averaging Process,
 1426 *Geoinfor Geostat: An Overview*, 1:2, doi:10.4172/gigs.1000103.
- 1427 Rohde, R. A. and Hausfather, Z. (2020), The Berkeley Earth Land/Ocean Temperature Record,
 1428 *Earth Syst. Sci. Data Discuss.*, <https://doi.org/10.5194/essd-2019-259>, in review
- 1429 Smith, T. M., R. W. Reynolds, T. C. Peterson, and J. Lawrimore (2008), Improvements to
 1430 NOAA's historical merged land–ocean surface temperatures analysis (1880–2006),
 1431 *Journal of Climate*, 21, 2283–2296, doi:10.1175/2007JCLI2100.1.
- 1432 Thorne, P. W., J. R. Lanzante, T. C. Peterson, D. J. Seidel and K. P. Shine (2011), Tropospheric
 1433 temperature trends: history of an ongoing controversy. *WIREs Clim Change*, 2: 66-88.
 1434 doi:10.1002/wcc.80
- 1435 Titchner, H. A., and N. A. Rayner (2014), The Met Office Hadley Centre sea ice and sea surface
 1436 temperature data set, version 2: 1. Sea ice concentrations, *J. Geophys. Res. Atmos.*, 119,
 1437 2864-2889, doi: 10.1002/2013JD020316.
- 1438 Yun, X., Huang, B., Cheng, J., Xu, W., Qiao, S., and Li, Q. (2019), A new merge of global
 1439 surface temperature datasets since the start of the 20th century, *Earth Syst. Sci. Data*, 11,
 1440 1629–1643, <https://doi.org/10.5194/essd-11-1629-2019>.
- 1441 Zhang, H-M, J. H. Lawrimore, B. Huang, M. J. Menne, X. Yin, A. Sánchez-Lugo, B. E. Gleason,
 1442 R. Vose, D. Arndt, J. J. Rennie, and C. N. Williams. (2019), Updated Temperature Data
 1443 Give a Sharper View of Climate Trends. *Eos*, 100,
 1444 <https://doi.org/10.1029/2019EO128229>
- 1445
- 1446

Figure 1.

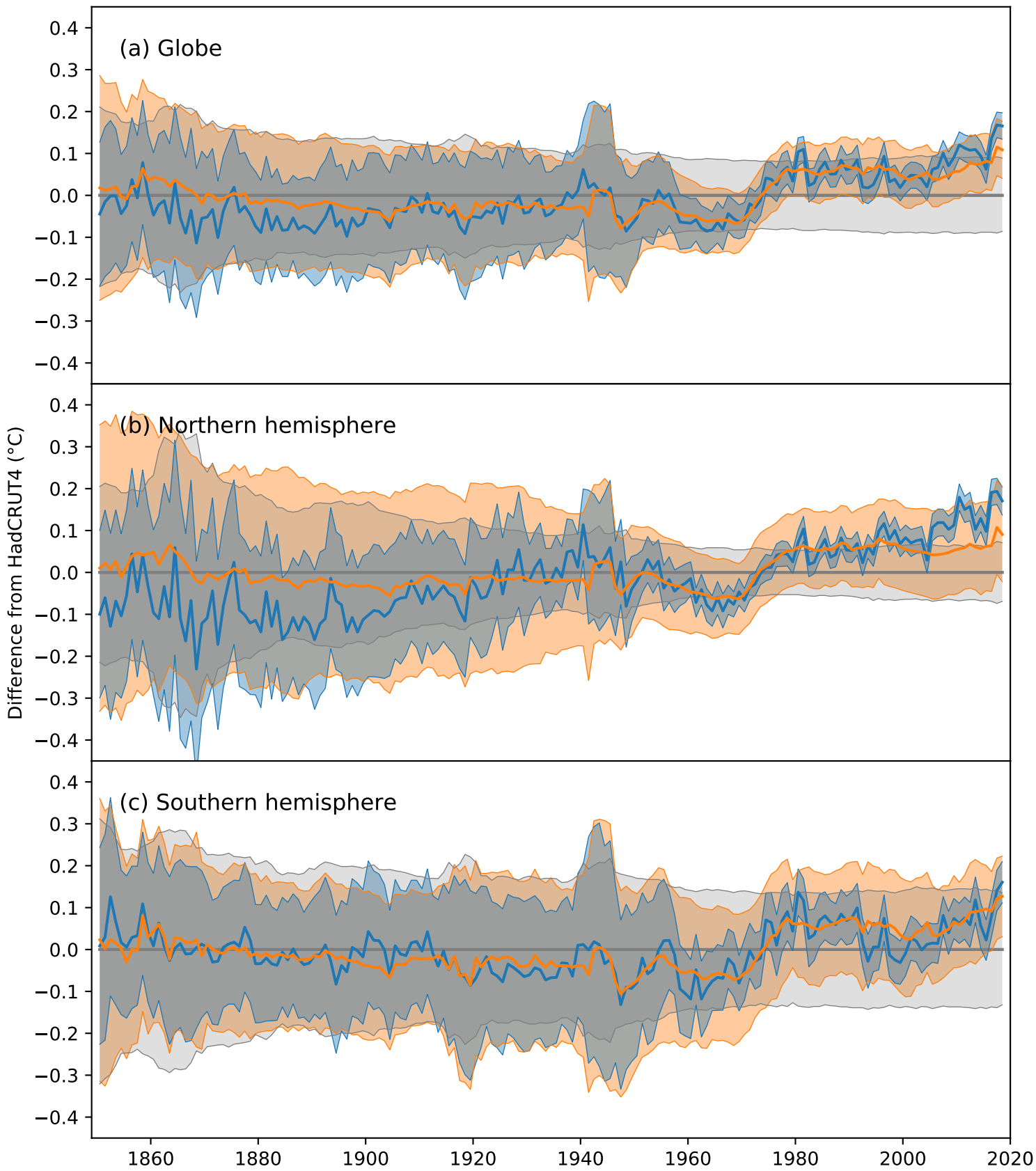


Figure 2.

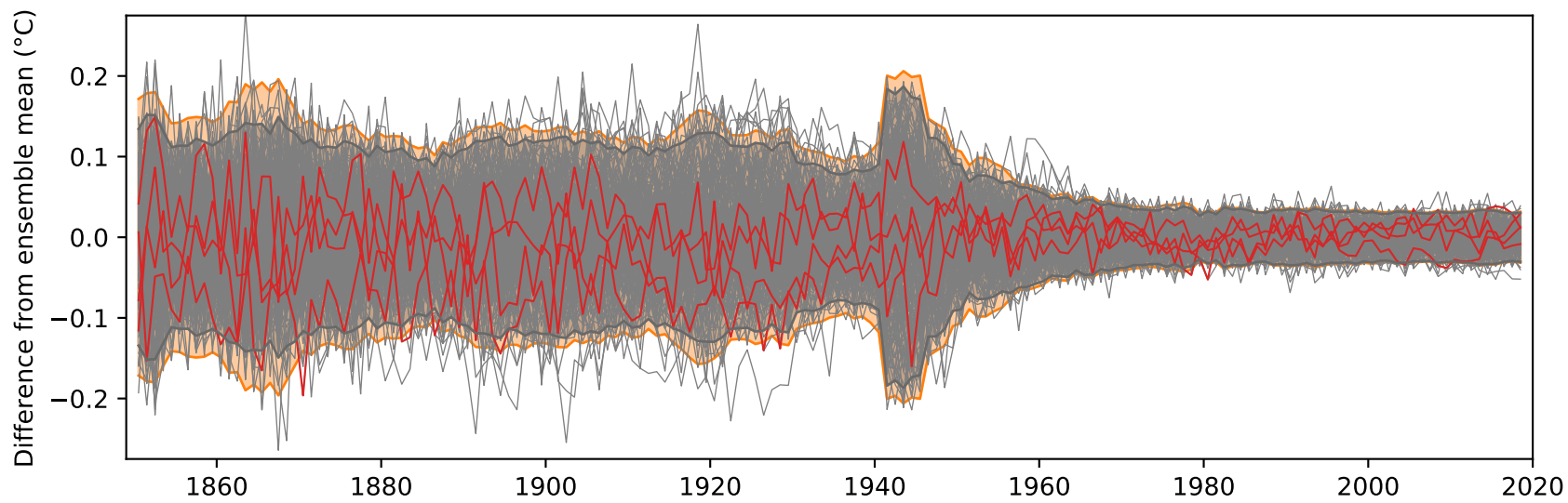
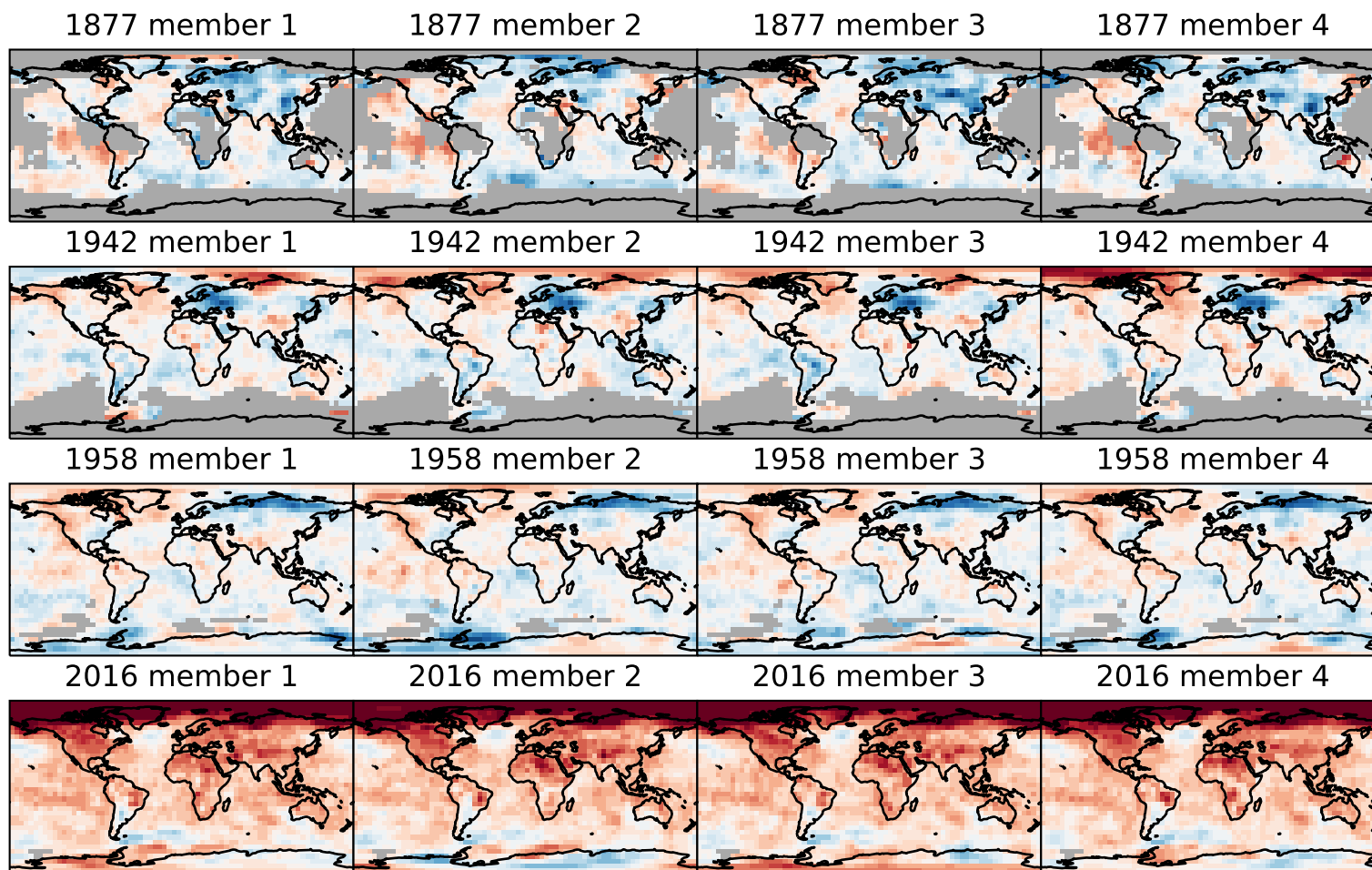


Figure 3.

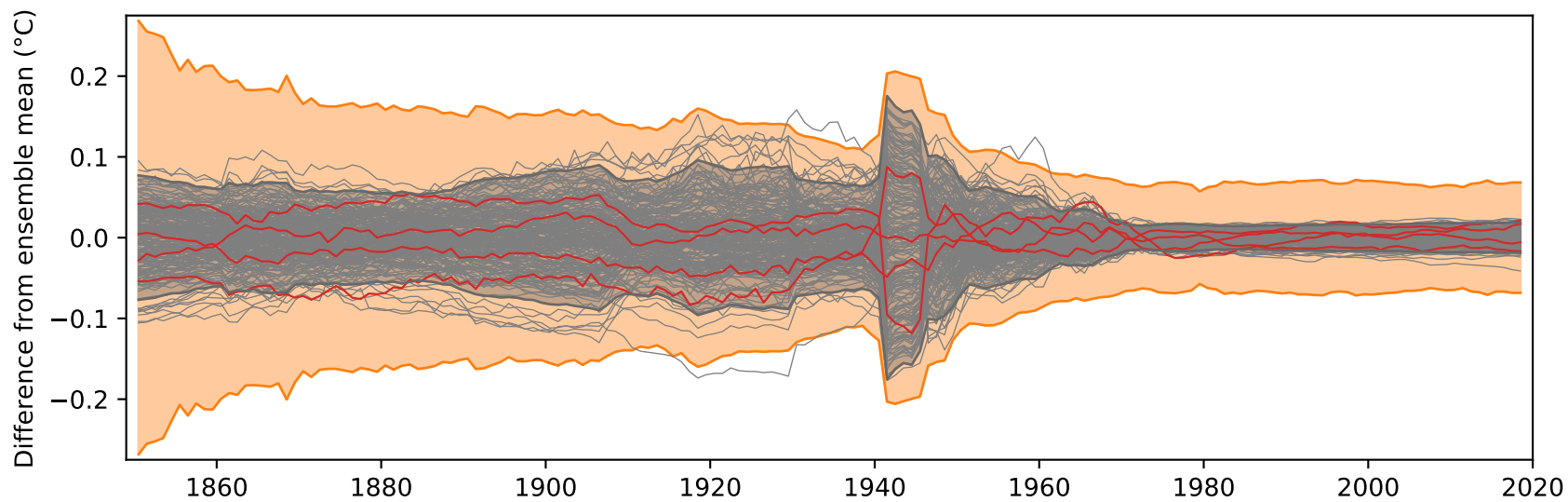
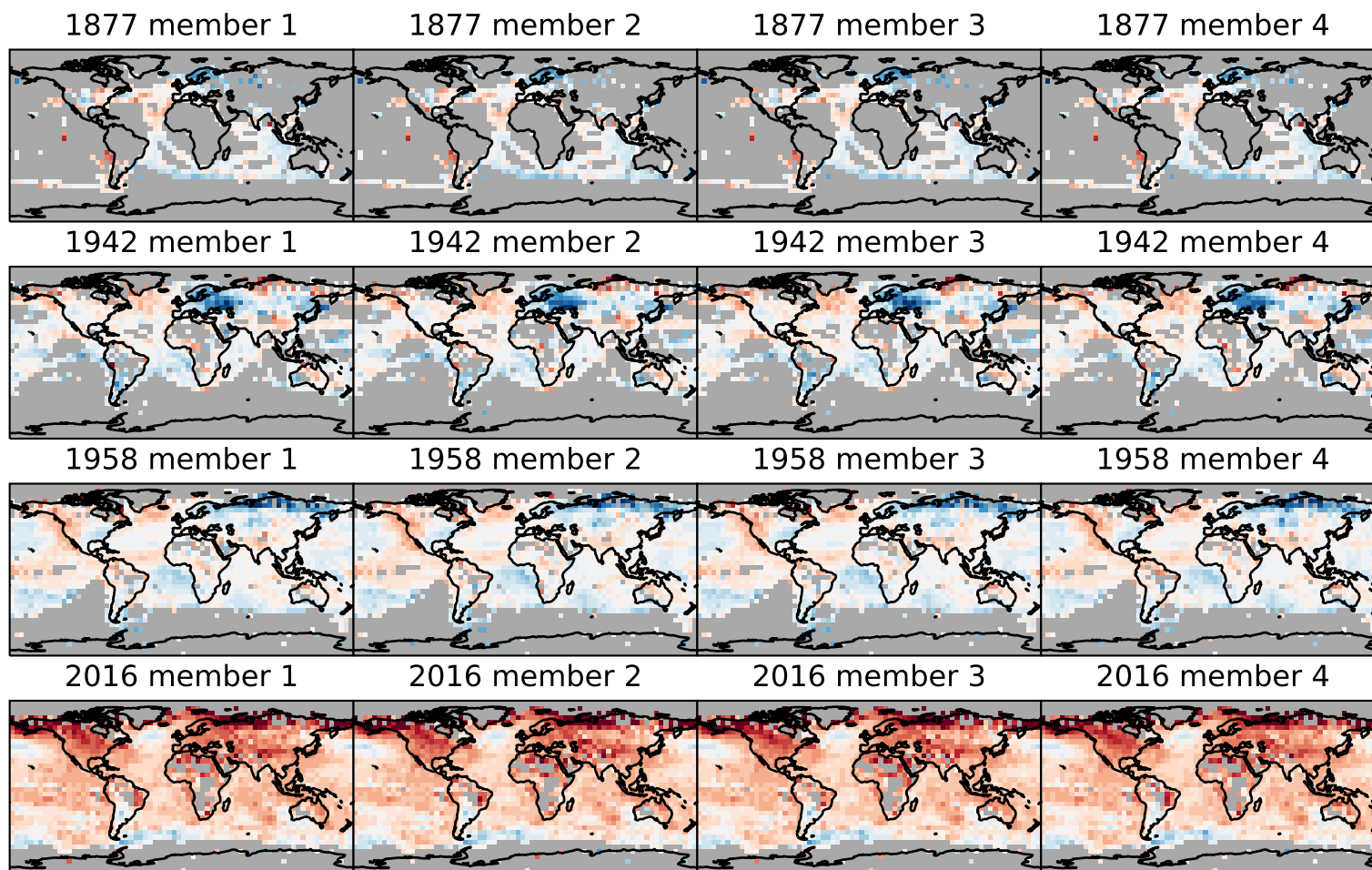


Figure 4.

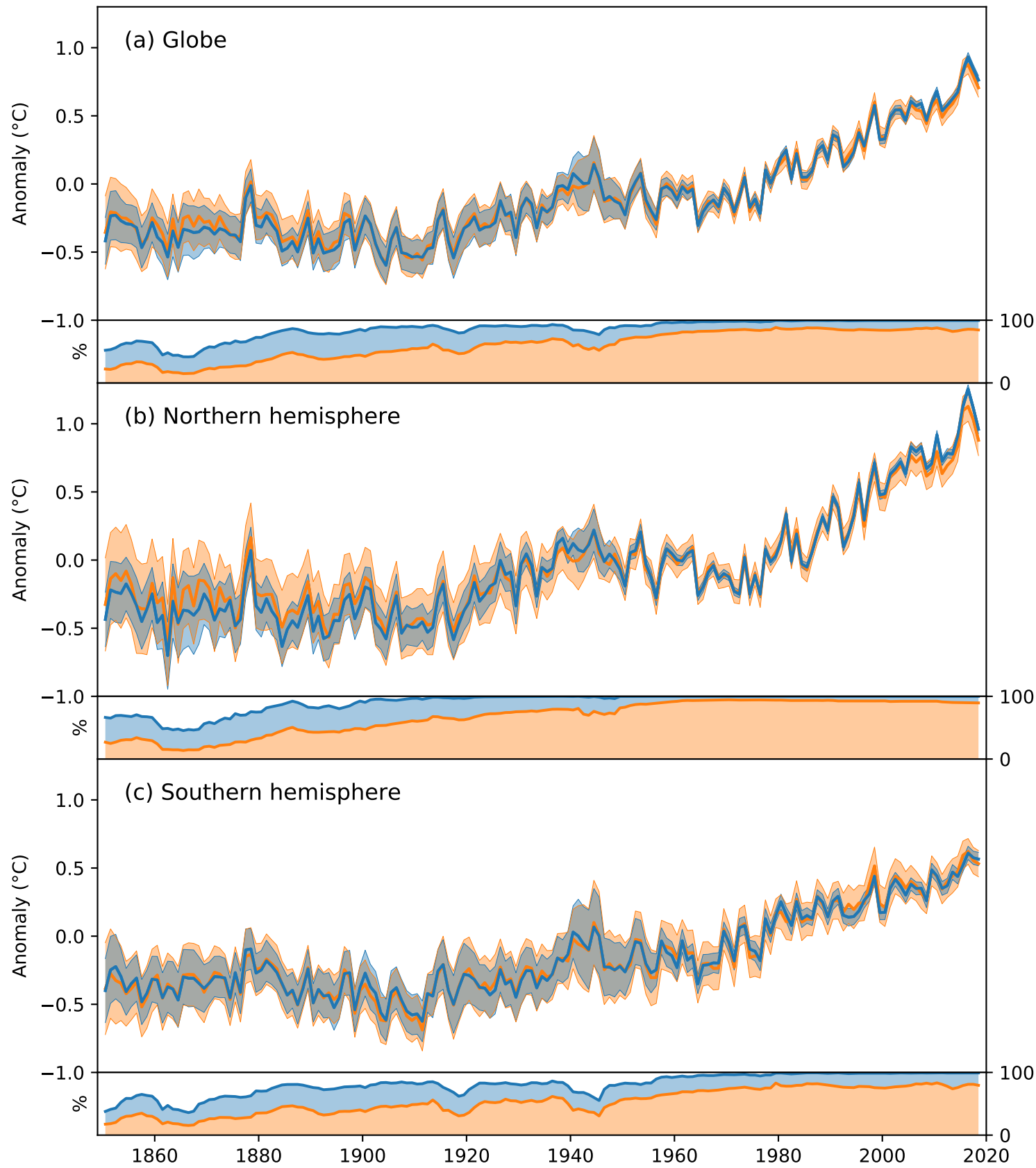


Figure 5.

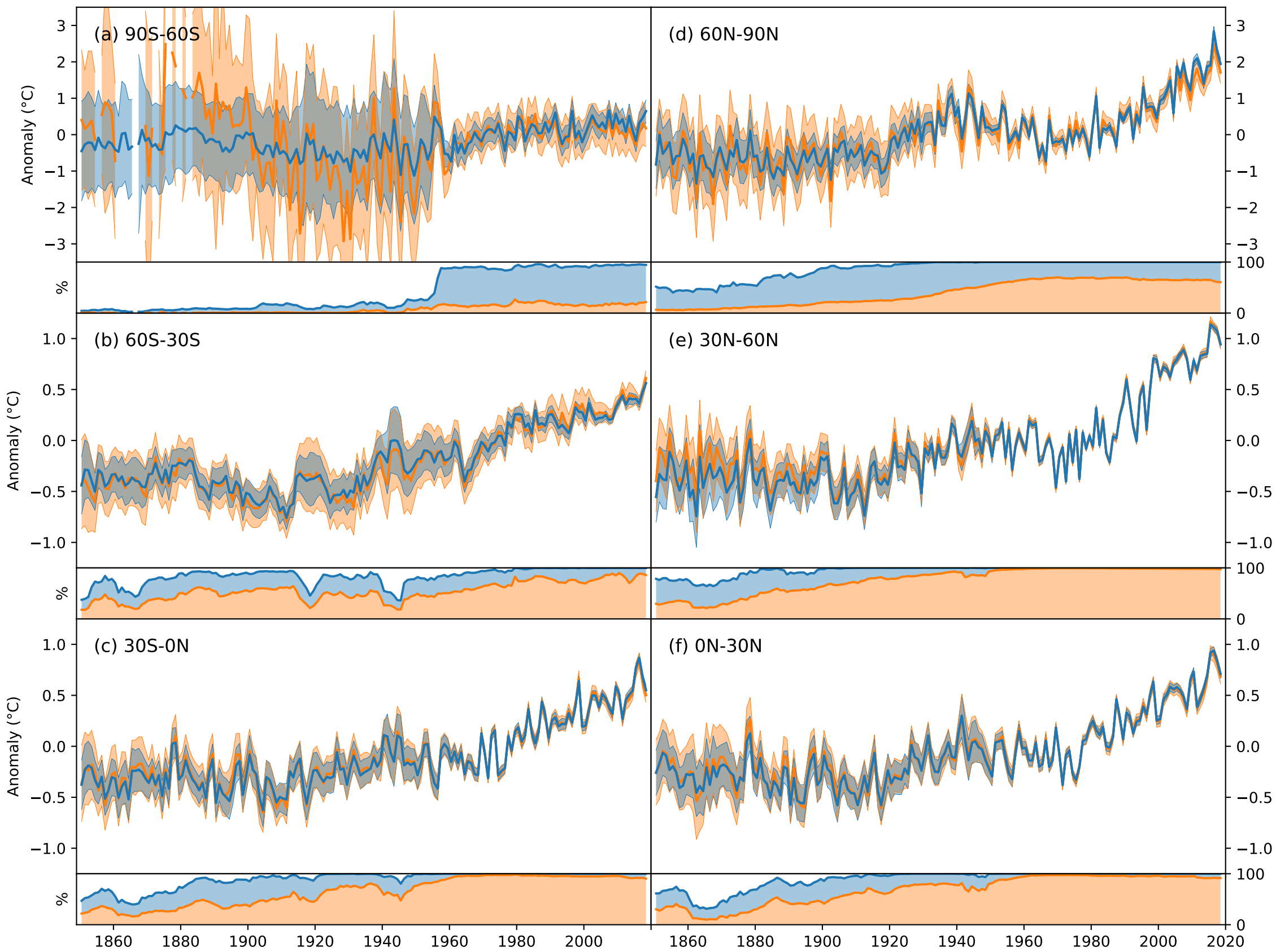


Figure 6.

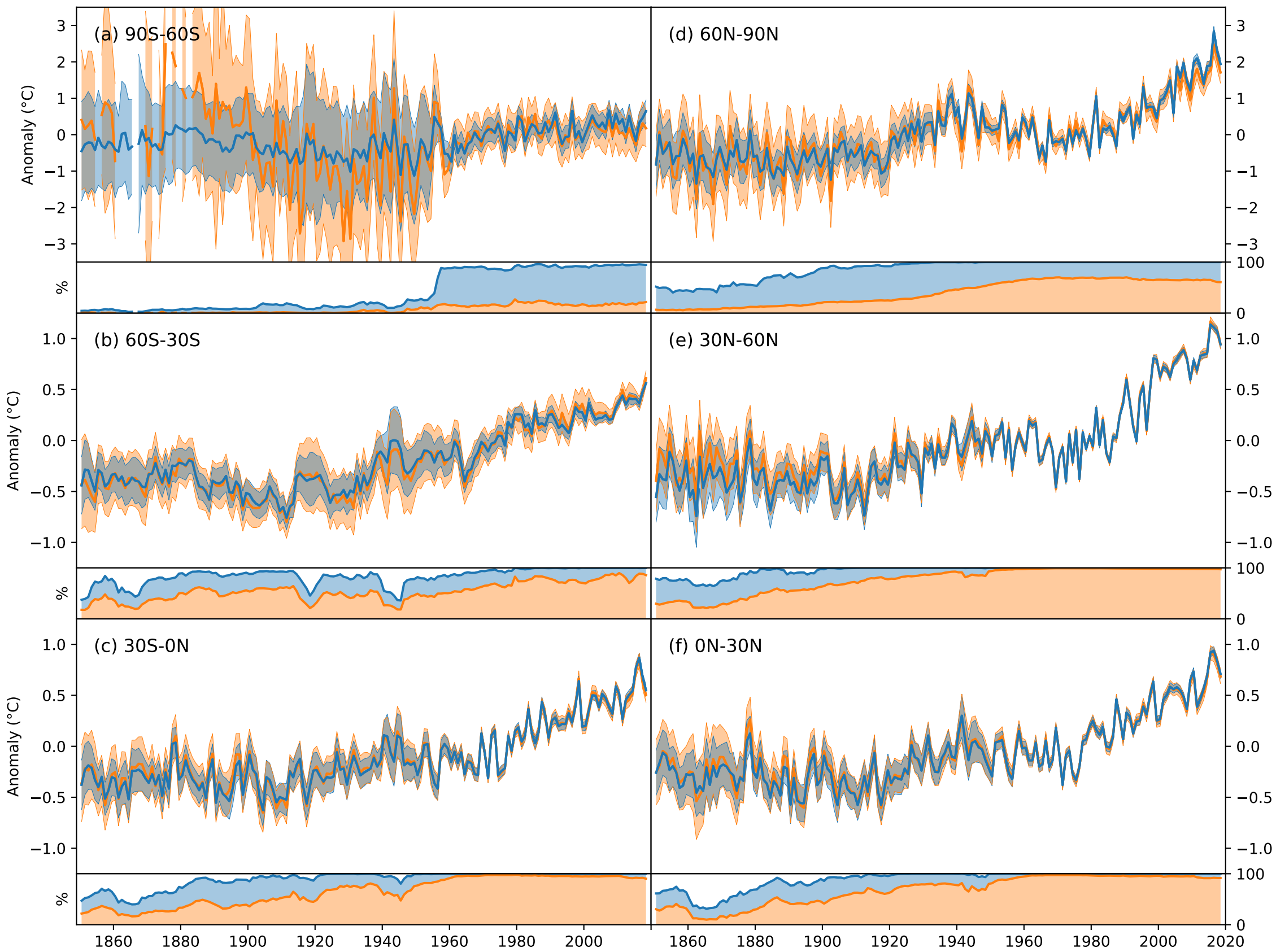


Figure 7.

Temperature anomaly ($^{\circ}\text{C}$ above 1961-1990 average)

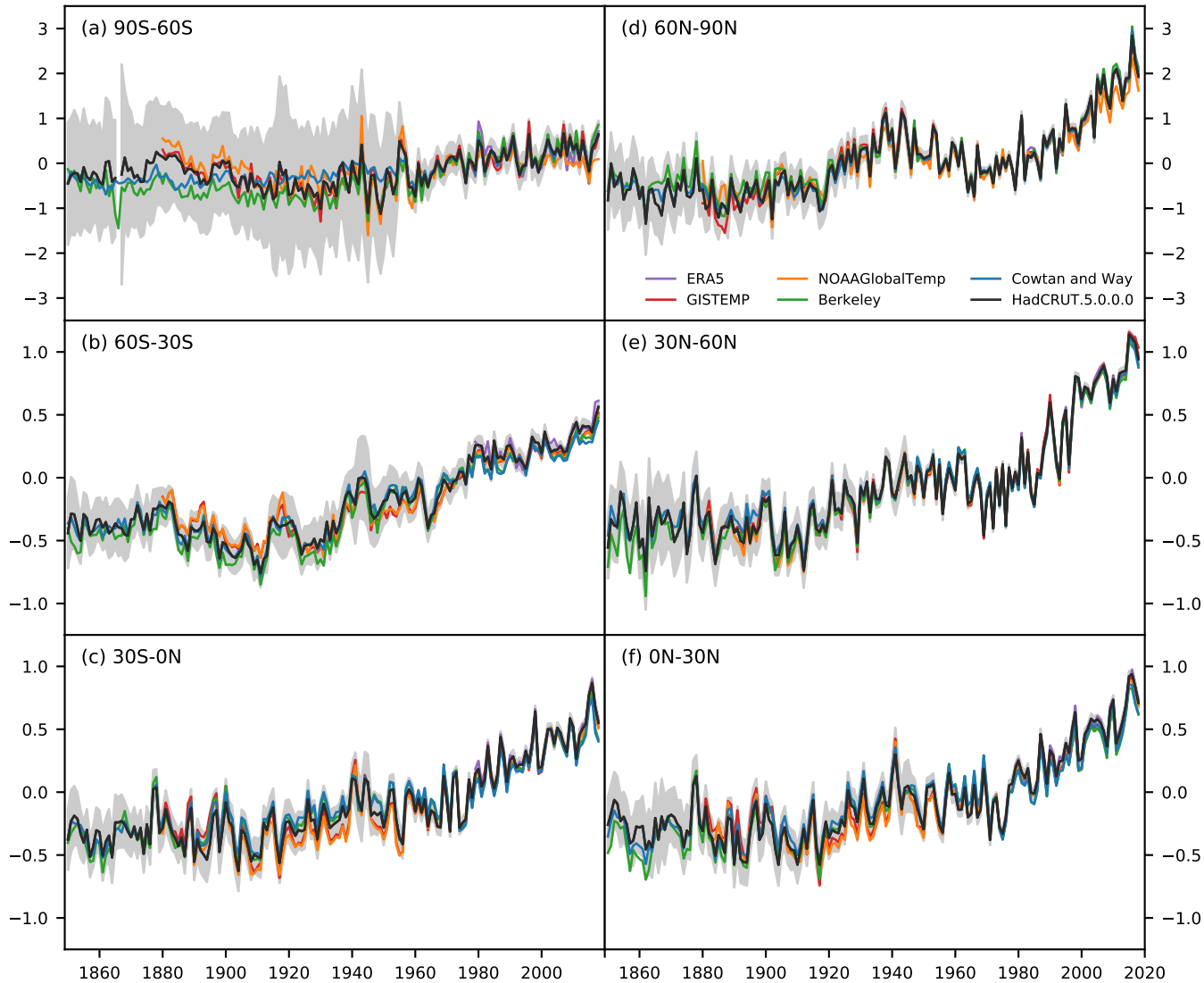


Figure 8.

