

1 **An ensemble data set of sea-surface temperature change from 1850: the Met Office**

2 **Hadley Centre HadSST.4.0.0.0 data set**

3 **Enter authors here: J.J. Kennedy¹, N.A. Rayner¹, C.P. Atkinson¹, R.E. Killick¹**

4 ¹Met Office Hadley Centre, Exeter, UK.

5 Corresponding author: John Kennedy (john.kennedy@metoffice.gov.uk)

6 **Key Points:**

- 7
- We describe the construction of HadSST.4.0.0.0, a climate data set of sea-surface

8 temperature change from 1850 to 2018.

 - A range of bias adjustments was generated to create an ensemble of SST data sets with

9 the ensemble spread partly constrained by oceanographic profile measurements.

 - New estimates reduce discrepancy between data sets during the mid 20th century and the

10 recent slowdown in warming, but highlight a divergence in the early 1990s.

11

12

13

14 **Abstract (max 250 words)**

15 One of the largest sources of uncertainty in estimates of global temperature change is that
16 associated with the correction of systematic errors in sea-surface temperature (SST)
17 measurements. Despite recent work to quantify and reduce these errors throughout the historical
18 record, differences between analyses remain larger than can be explained by the estimated
19 uncertainties.

20 We revisited the method used to estimate systematic errors and their uncertainties in version 3 of
21 the Met Office Hadley Centre SST data set, HadSST. Using comparisons with oceanographic
22 temperature-profiles, we make estimates of biases associated with engine-room measurements
23 and insulated buckets and constrain the ranges of two of the more uncertain parameters in the
24 bias estimation: the timing of the transition from uninsulated to insulated buckets in the mid-20th
25 century and the estimated fractions of different measurement methods used.

26 Here, we present HadSST.4.0.0.0, based on release 3.0.0 and 3.0.1 of the International
27 Comprehensive Ocean-Atmosphere Data Set supplemented by drifting buoy measurements from
28 the Copernicus Marine Environmental Monitoring Service. HadSST.4.0.0.0 comprises a 200-
29 member “ensemble” in which uncertain parameters in the SST bias scheme are varied to generate
30 a range of adjustments. The evolution of global average SST in the new data set is similar to that
31 in other SST data sets and the difference between data sets is reduced during the mid-20th
32 century. However, the changes also highlight a discrepancy in the global-average difference
33 between adjusted SST and marine air temperature in the early 1990s and hence between
34 HadSST.4.0.0.0 and, the NOAA SST data set, ERSSTv5.

35

36 **1 Introduction**

37 Sea-surface temperature (SST) has been designated an essential climate variable (ECV, Bojinski
38 et al. 2014) that “critically contributes to the characterization of Earth’s climate” (WMO
39 [https://public.wmo.int/en/programmes/global-climate-observing-system/essential-climate-](https://public.wmo.int/en/programmes/global-climate-observing-system/essential-climate-variables)
40 [variables](https://public.wmo.int/en/programmes/global-climate-observing-system/essential-climate-variables)). The Global Observing Systems Information Centre (GOSIC) website states,
41 *“Together with air temperature over land, sea-surface temperature is the most important*
42 *variable for determining the state of the climate system”*. It is a key variable for detection of
43 climate change and assessing the relative importance of anthropogenic and natural influences on
44 global climate (<http://www.gosic.org/content/gcos-oceanic-surface-ecv-sea-surface-temperature>).
45 Atmospheric and oceanic reanalyses (e.g. Kalnay et al. 1996, Hersbach et al. 2015, Carton and
46 Geise 2008, Compo et al. 2011), which are some of the most widely used and cited tools in
47 weather and climate science, typically use SST data sets to provide a lower (or upper) boundary
48 condition. Consequently, there is some value in understanding the long-term evolution of SST
49 and its uncertainties.

50
51 Historical SST measurements are to be found, digitized, in great numbers in the International
52 Comprehensive Ocean-Atmosphere Data Set (ICOADS) alongside many other marine
53 meteorological variables. The two most recent major releases of ICOADS are release 2.5
54 (Woodruff et al. 2010), which contains 261 million records and covers 1662-2007 and release
55 3.0 (Freeman et al. 2016), which contains over 455 million individual marine reports and covers
56 1662-2014. Although SST measurements are few in the very early record, they become much
57 more numerous in the latter half of the nineteenth century. Thus far, only ERSSTv5 (currently
58 1880-2018, Huang et al. 2017) has made use of the far greater number of measurements

59 available in ICOADS release 3.0. Other major historical gridded SST data sets, which run from
60 the latter half of the 19th century, such as COBE-SST-2 (1880-2010, Hirahara et al. 2014),
61 ERSSTv4 (currently 1880-2018, Huang et al. 2015, Liu et al. 2015), and HadSST.3.1.1.0
62 (currently 1850-2018, Kennedy et al. 2011a; Kennedy et al. 2011b) use ICOADS release 2.5.

63
64 One particular difficulty associated with making long data sets based on SST measurements for
65 use in climate analyses is that the technology used to measure SST has changed so much over the
66 past one and a half centuries (Kent et al. 2010, Kent et al. 2017). Even subtle changes in the way
67 that measurements are carried out can lead to systematic errors in the measured trends and the
68 historical changes have not been especially subtle. The magnitude of the estimated errors are of
69 order 0.1-1.0°C, similar to climatic variations over the same period (Hartmann et al. 2013). It is
70 therefore necessary to correct these systematic errors and quantify the residual uncertainties to
71 better understand what actually happened. Each of the previously mentioned data sets – COBE-
72 SST-2, ERSSTv4/v5, and HadSST.3.1.1.0 - applies adjustments to correct systematic errors in
73 the data and provides some estimate of the uncertainty.

74
75 Kennedy et al. (2011b) generated an ensemble of one hundred members which comprise the
76 HadSST.3.1.1.0 data set. They calculated a range of corrections by varying poorly-constrained
77 parameters in their bias-adjustment scheme. They used metadata from a number of sources (for
78 example, instructions to marine observers) to assign a measurement method to each observation
79 and took estimates of the systematic errors associated with each measurement method from the
80 literature. The residual uncertainty was combined with uncertainties from other sources of error
81 such as sampling and local measurement errors (Kennedy et al. 2011a and b). An ensemble

82 approach to quantifying uncertainties was also used for the ERSSTv4 data set (Liu et al. 2015,
83 Huang et al. 2016). The parameters they varied to generate the ensemble were associated with all
84 steps in the data-set creation process and only a subset relates directly to the correction of
85 systematic errors. The ERSST SST adjustments are based on comparisons with gridded Night
86 Marine Air Temperature (NMAT) data from the HadNMAT2 data set (Kent et al. 2013) and on
87 estimated differences between SSTs measured by ships and buoys. Hirahara et al. (2014) did not
88 use an ensemble approach, but did provide statistical estimates of the uncertainties in the
89 adjustments applied to the COBE-SST-2 data set. In addition, they used the data to improve the
90 estimates of some of the uncertain parameters identified in Kennedy et al. (2011b).

91

92 Despite the efforts of these researchers, significant differences remain between the data sets even
93 at a global scale. The most notable differences (Kent et al. 2017) are between HadSST.3.1.1.0
94 and ERSSTv5 in the period around the Second World War and since the late 1990s. Larger
95 differences earlier in the record are within the joint uncertainty range.

96

97 The period around the Second World War is a period marked by profound uncertainty in the SST
98 record. The war years saw a prolonged El Niño (from late 1939 to 1942, Brönnimann 2005)
99 bringing a peak in global temperatures. The end of the war coincided with a shift in the phase of
100 the Pacific Decadal Oscillation (Newman et al. 2016) with an ensuing period of relatively stable
101 or declining global temperature. A change in both the pattern of international shipping and the
102 composition of data sources available in ICOADS (Thompson et al. 2008) occurs at the same
103 time, confounding a straightforward understanding of the events. The following decades, from
104 1950 to 1970, witnessed large and poorly documented changes in the way that measurements

105 were made with the development of high-tech insulated buckets and a long-term shift towards
106 measurements being made in the engine rooms of ships: also known as ERI measurements. The
107 “I” has been variously taken to mean intake, inlet and injection.

108

109 Even prior to the Second World War, when buckets were the primary means of sampling
110 seawater to measure SST on ships, there are geographical and seasonal differences between the
111 adjustments in ERSSTv5 and HadSST.3.1.1.0. The differences arise from the assumed
112 dependence of the biases on weather conditions. Pre war, ERSST uses adjustments that depend
113 on the air-sea temperature difference alone. The adjustments used in HadSST.3.1.1.0 and COBE-
114 SST-2 (Folland and Parker 1995, Rayner et al. 2006), which assume evaporative cooling from
115 the wet surfaces of the bucket, depend not only on the air-sea temperature difference, but also on
116 solar radiation and, critically, the wet-bulb depression (Carella et al. 2017b). The Folland and
117 Parker (1995) model has recently been assessed in the laboratory by Carella et al. (2017b) who
118 found that the model performed well when conditions were known and controlled, but noted that
119 measurement conditions on board ship were typically neither of these things.

120

121 Differences between data sets in the modern period, marked by a slow dwindling of the
122 Voluntary Observing Ship (VOS) fleet and the widespread deployment of drifting buoys, are
123 largely within the joint uncertainties of the various SST data sets. However, the period from
124 2000 to 2013, during which global temperatures increased at a lower rate than some expected,
125 has been intensively studied (Medhaug et al. 2017) and high demands for accuracy have been
126 made of the SST data sets. It would be advantageous to have a more reliable estimate of SSTs

127 during this well-observed period (Karl et al. 2015, Hausfather et al. 2017) the better to establish
128 exactly what happened.

129

130 No part of the SST record is simple to understand or without some little mystery of its own.
131 Therefore, the aims of this paper are to revisit and improve the adjustments applied to
132 HadSST.3.1.1.0 and explore the residual uncertainties, paying particular attention to how the
133 measurements were made. We start by describing the data sources used in the analysis in Section
134 2. Section 3 explains how the data are aggregated onto a regular grid and how uncertainties
135 associated with the uncorrelated measurement errors and under-sampling are estimated. The bias
136 adjustments, which deal with other, correlated measurement errors and the creation of the
137 ensemble, are the focus of Section 4 with some of the technical details included in the Appendix.
138 Results are in Section 5 before we finish up in Section 6 with the presentation of the new
139 HadSST.4.0.0.0 data set, some discussion and general conclusions.

140

141 Throughout this paper, we frequently cite Kennedy et al. (2011a, 2011b and 2011c) as well as
142 Rayner et al. 2006 as this data set builds directly on these papers. We refer to these as K11a,
143 K11b, K11c and R06 for brevity. Also frequently cited are Folland and Parker (1995), hereafter
144 FP95, and Smith and Reynolds (2002), which we shorten to SR02.

145 **2 Data**

146 We use various data sets in the analysis and for comparison and validation; they are described in
147 the following subsections. The main analysis is based on the International Comprehensive
148 Ocean-Atmosphere Data Set (ICOADS, Freeman et al. 2017) and water temperatures from
149 HadIOD.1.2.0.0, the Met Office Hadley Centre Integrated Ocean Database (Atkinson et al.

150 2014). For inter comparison, we use ERSSTv4 (Huang et al. 2015), ERSSTv5 (Huang et al.
151 2017) and COBE-SST-2 (Hirahara et al. 2014) which have already been introduced. We also use
152 HadNMAT.2.0.1.0 (Kent et al. 2013). For validation we use independent satellite SST retrievals
153 from the ATSR (Along Track Scanning Radiometer) Reprocessing for Climate (ARC) data set
154 (Merchant et al. 2012) as well as some instrumentally homogeneous (Hausfather et al. 2017)
155 subsets of the HadIOD and ICOADS data sets. First, however, it is useful to determine exactly
156 what it is we mean by sea-surface temperature.

157 **2.1 Which sea-surface temperature?**

158 Traditionally, long-term *in situ* SST data sets have been considered to be representative of a
159 loosely defined “bulk” SST, which covers a range of measurements made in the upper 10m or so
160 of the water column. However, the daily formation and erosion of a stably-stratified near-surface
161 warm layer in the oceans, particularly during calm, sunny conditions, can lead to strong
162 temperature gradients in the upper 10m (Kawai and Wada 2007) and make it harder to reconcile
163 measurements of water temperature made at different depths. Satellite retrievals of SST are
164 especially prone to this as they are sensitive to temperatures in a very shallow layer where
165 diurnal warming is most pronounced. For instruments measuring in the infra red, this layer is
166 measured in micrometres and its temperature is referred to as the “skin” temperature.

167 Uncertainties associated with modern measurement systems are now sufficiently small that
168 temperature variations with depth are readily detectable in the aggregate and need to be
169 accounted for where detailed comparisons are made (see e.g. Merchant et al. 2012).

170

171 Donlon et al. (2007) recommend that all SST measurements be accompanied by an estimate of
172 the depth at which the measurement was made. However, this information is rarely available for

173 historical observations and, where it is available, it does not usually take into account changes
174 due to variable loading of the ship. Even drifting buoys from a single manufacturer, where the
175 design and materials are identical, will measure at a varying depth owing to continual movement
176 of the water and the potential loss of the drogue.

177

178 In this paper, rather than specify that SSTs are estimated for a particular depth, we will instead
179 use SST measurements from drifting buoys as our reference. This is conventionally reckoned
180 equivalent to an SST measurement at an approximate depth of 20cm. We make use of near-
181 surface water temperatures measured at a range of depths, as described in the following
182 subsections, but the aim throughout is to adjust these measurements so that they would closely
183 match coincident observations from drifting buoys even if they occurred long before drifting
184 buoys were first dreamed of.

185

186 A final note on nomenclature. Throughout this paper we use the words “error” and “uncertainty”
187 as they are defined in Annex B of the Guide to the Expression of Uncertainty in Measurement
188 (JCGM 2008).

189

190 **uncertainty** (of measurement): *parameter, associated with the result of a measurement, that*
191 *characterizes the dispersion of the values that could reasonably be attributed to the measurand.*

192

193 **error** (of measurement): *result of a measurement minus a true value of the measurand.*

194

195 A *measurand* is a “particular quantity subject to measurement”.

196

197 **2.2 Surface meteorological data**

198 The data and metadata used in this analysis are from Release 3.0 of the International
199 Comprehensive Ocean-Atmosphere Data Set (ICOADS) for the period 1850-2014 (Freeman et
200 al. 2017; downloaded version Research Data Archive 2016). An update to 2017 uses the near-
201 real time data from ICOADS release 3.0.1. We will use ICOADS release 3.0.1 as the basis for
202 monthly updates of the data set.

203

204 Due to a drop off in drifting buoy observations in ICOADS release 3.0.1 which followed the
205 switch of data transmission codes from TAC (Traditional Alphanumeric Codes) to BUFR
206 (Binary Universal Form for the Representation of meteorological data) in November 2016, we
207 supplement the ICOADS data from January 2016 onwards with near-real-time drifting buoy
208 observations downloaded from CMEMS (Copernicus Marine Environmental Monitoring
209 Service, [ftp://nrt.cmems-
210 du.eu/Core/INSITU_GLO_NRT_OBSERVATIONS_013_030/monthly/drifter/](ftp://nrt.cmems-du.eu/Core/INSITU_GLO_NRT_OBSERVATIONS_013_030/monthly/drifter/)). The initial
211 download was made on 5 April 2018 and then regular downloads are made each month to gather
212 data for the preceding month. The drifting buoys from this near-real-time source completely
213 replace the drifting buoys from ICOADS release 3.0.1 in our analysis in the overlapping months.
214 This increases the data volume and observational coverage significantly from November 2016.

215

216 We filtered the data to remove coastal stations, non-standard moored buoys (principally around
217 the coast of the US) and other non-standard platforms like oil rigs and research stations. These
218 sources cover a relatively small area and vary widely in design making the biases heterogeneous

219 and difficult to assess. C-MAN and many US Coastal moored buoys were excluded because, as
220 well as being of diverse type, they are often found in estuaries and lagoons that are not
221 representative of a wider area and certainly not of the areas typically sampled during the
222 climatological base period 1961 to 1990. Many coastal moored buoys also produce very high
223 data volumes – sometimes making several measurements an hour – and can potentially skew the
224 processing. A list of the excluded ICOADS platform IDs is available along with the
225 HadSST.4.0.0.0 data set.

226

227 Oceanographic measurements from the World Ocean Database (WOD) were also removed by
228 excluding ICOADS deck 780 (1850-2014). Excluding WOD measurements from the processing
229 means that we can use sub-surface measurements as an independent data set for assessing biases
230 and for validation (e.g. Gouretski et al. 2012, Huang et al. 2018). The remaining data were
231 quality controlled (QC'd) to remove outliers and low-quality measurements (an update of R06,
232 <https://github.com/ET-NCMP/MarineQC>). Figure 1 shows the number of observations passing
233 QC for each month from January 1850 to December 2018.

234

235 R06 describes the creation of the climatology we use to calculate the gridded anomalies. It is
236 based on *in situ* measurements made between 1961 and 1990. The climatology has a resolution
237 of 1° of latitude, 1° of longitude and 5 days. A standard 5-day period is a pentad. The first pentad
238 of each year is 1-5 January. The calendar is divided into pseudo months. Each pseudo-month has
239 six pentads except August, which has seven (hence the annual peaks in Figure 1). Leap days are
240 accommodated by extending the pentad in which they fall.

241 **2.2.1 Initial metadata assignment**

242 We assign a measurement method to each observation in ICOADS. Any particular report could
243 be either a buoy measurement, a bucket measurement, an ERI measurement, a measurement
244 made with a hull sensor or else unknown. Where we could not definitively assign a single
245 method, a fractional assignment was attempted based on the fraction of the recruiting country's
246 fleet that used each method (see also K11b). Fractions were estimated based on ships recorded in
247 WMO Publication 47 (WMO Pub. 47, Kent et al. 2007) for that year. An assignment from the
248 ICOADS metadata (SI or SIM indicating bucket, ERI or hull sensor) was preferred. Fractional
249 assignments are always incorrect at the level of individual reports, but should give representative
250 averages when aggregating large numbers of observations. The assignments are uncertain, even
251 where there is a definitive assignment, and we refine the estimates using comparisons between
252 the ship data and oceanographic profile data in Section 4.1.4.

253

254 The procedure for assigning metadata to a particular ICOADS report is as follows. The
255 procedure terminates as soon as an assignment is made (abbreviations in brackets refer to the
256 variable names in the IMMA, International Maritime Meteorological Archive format,
257 documentation for ICOADS release 3.0.0 <http://icoads.noaa.gov/e-doc/imma/R3.0-imma1.pdf>):

258

- 259 1. If the ICOADS platform type (PT) was 6 or 7, we assign the observation to be a moored
260 or drifting buoy measurement respectively.
- 261 2. We assigned US ships (C1 = 2) from Deck (DCK) 128 to be ERI in 1968, 1969, 1972 and
262 1973.

- 263 3. UK Royal Navy data, Deck 245, and Russian data from Deck 732 (following Carella et
264 al. 2018) were assigned to be ERI.
- 265 4. If an SST measurement method (SI) was present in ICOADS and indicated a bucket, ERI
266 or hull measurement, we used it.
- 267 5. If an SST measurement method (SIM) was present in the ICOADS metadata attachment
268 and indicated a bucket, ERI or hull measurement, we used it.
- 269 6. Between 1939 and 1945, we set all reports that had not been assigned a measurement
270 method in steps 1-5 to unknown.
- 271 7. Before 1939, all reports that had not been assigned a measurement method in steps 1-6
272 were set to bucket.
- 273 8. If a recruiting country (C1) was present in the ICOADS attachment, we used it to assign
274 weights to each report for bucket, ERI, hull and unknown according to the fraction of
275 ships that took each type of measurement in WMO Pub. 47 for that country and year. We
276 assumed that US ships with unknown measurement method were ERI. If WMO Pub. 47
277 was not available for that year, we used the next available year after the year of interest.
- 278 9. If the deck (DCK) could be linked to reports from a particular country (see K11b for
279 details) between 1956 and 1996, then we used it to assign weights to each report of
280 bucket, ERI, hull and unknown according to the fraction of ships that took each type of
281 measurement in WMO Pub. 47 for that country and year. We assumed that US ships with
282 unknown measurement method were ERI. If WMO Pub. 47 was not available for that
283 year, the next available year was used.
- 284 10. Any report not assigned a measurement method in steps 1-9 we set to unknown.

285 At the end of this process, each SST measurement has an assigned measurement method, or
286 fractional assignment. The assignments are provisional and uncertain. In some cases,
287 observations from a ship will be listed as being made with a bucket, when they were in fact ERI
288 measurements. In section 4.1.4, we attempt to estimate what fraction of bucket measurements are
289 mis-identified in this way.

290 **2.3 Near-surface sub-surface measurements from HadIOD**

291 HadIOD.1.2.0.0 (Atkinson et al. 2014) is an integrated database of temperature and salinity
292 measurements from oceanographic sources at various depths (from the EN version 4.2.0 data set;
293 Good et al. 2013) combined with surface observations from ICOADS release 2.5.1 (Woodruff et
294 al. 2011). Oceanographic measurements in HadIOD were made with a variety of instruments
295 including Conductivity-Temperature-Depth (CTD), eXpendable BathyThermograph (XBT),
296 Mechanical BathyThermograph (MBT) and Argo devices. For a review of the characteristics of
297 these devices, see Abraham et al. (2013). Each observation has an overall quality flag, an
298 estimated uncertainty and a bias adjustment or bias adjustments. We only use the oceanographic
299 profile measurements from HadIOD between 1930 and 2018 and of these we use only those that
300 were made in the upper 10m of the water column. Measurements in the upper 10m can provide a
301 reasonable approximation for the sea-surface temperature (see Section 2.1). We set aside
302 measurements from Argo floats so that we can use them for independent validation (see Section
303 5.2 and Appendix 3.1).

304

305 A number of XBT and MBT adjustments are included in the HadIOD database. The adjustments
306 are required to correct for known biases in these measurements independent of the SST biases
307 we explore in this paper. We used MBT and XBT adjustments from four analyses (Gouretski and

308 Reseghetti 2010, Gouretski 2012, Cowley et al. 2013 and Levitus et al. 2009). Two of the
309 analyses provide estimates of both MBT and XBT adjustments (Levitus et al. 2009 and
310 Gouretski and Reseghetti 2010) and we use these as given. Two analyses provide only XBT
311 adjustments (Gouretski 2012 and Cowley et al. 2013). We combined the two XBT-only analyses
312 with each of the two MBT adjustments from the combined analyses to provide four new XBT-
313 MBT combinations, bringing the total number of sets of corrections to six. The six sets of
314 adjusted XBT-MBT data together with measurements from Conductivity-Temperature-Depth
315 (CTD) instruments and bottles are gridded as for the SSTs (Section 3) to make six near-surface
316 reference data sets that we use to provide a set of baselines – albeit uncertain baselines – for the
317 SST bias assessment. We assume that, once corrected, the oceanographic data are less biased
318 than the SST measurements are. However, we note that the adjustment of oceanographic data is
319 also an area of ongoing research with uncertainties all of its own (Abraham et al. 2013, Cheng et
320 al. 2016).

321

322 **2.4 Marine air temperature**

323 The air temperature above the oceans is physically related to the underlying sea-surface
324 temperature. Consequently, sea-surface temperature changes are often used as a proxy for marine
325 air temperature (MAT) changes, for example, in the calculation of global average temperature
326 anomalies (e.g. Morice et al. 2012). Huang et al. (2015) argued, based on the behavior of a
327 particular climate model (the GFDL coupled model, CM2.1), that the difference between MAT
328 and SST anomalies at a global scale was more or less constant, changing by less than 0.1°C in a
329 century. However, other authors have noted that trends in MAT and SST anomalies can diverge,
330 albeit by a relatively small amount on multidecadal time scales (Cowtan et al. 2015). At smaller

331 scales, MAT and SST anomalies can diverge by tens of degrees particularly close to land and sea
332 ice.

333

334 Nonetheless, over longer periods and at larger scales, changes in MAT and SST are closely
335 related. Indeed, SR02, FP95 and other related papers used this relationship to adjust for biases in
336 the SST record. However, MAT measurements are not without problems of their own. The
337 principal problem is solar heating of the ship, which biases MAT measurements during the day
338 (Berry et al. 2004) and is typically solved by restricting the use of MAT measurements to those
339 made at night: so called night marine air temperature (NMAT) measurements. The increasing
340 size of ships and the height of temperature sensors above the sea-surface has led to a creeping
341 cold bias in MAT measurements. Warm biases have also been detected during the 19th century
342 and the Second World War, and are thought to be caused by non-standard sensor exposure; for
343 example, reading the thermometer inside rather than on deck. In the creation of the HadNMAT2
344 data set (1880-2010), Kent et al. (2013) adjusted the data, or excluded certain periods, regions,
345 and data subsets, to account for these biases.

346

347 **2.5 Instrumentally homogeneous data sets**

348 Hausfather et al. (2017) used the term “instrumentally homogeneous” data sets to describe SST
349 data sets that are based on a single type of instrument or group of closely-related instruments,
350 which they considered to be more homogeneous than the general *in situ* SST record. They used
351 these data sets to assess the stability of global-average SST records over the period 1995-2016.

352

353 Records that consist of a single type of instrument minimize artificial drifts, shifts or jumps
354 caused by changes in instrumentation. Ideally, measurements from the instruments should also
355 be of high quality, with demonstrably good accuracy and stability. Finally, the measurements
356 should be as independent as possible from the record they are being used to assess. We use three
357 instrumentally homogeneous records based on: Argo floats, Along-Track Scanning Radiometer
358 SST retrievals and buoys. We describe each of these in turn.

359

360 **2.5.1 Argo**

361 Argo floats are autonomous profiling floats, which move with the prevailing currents at a typical
362 “parking depth” of 1000m descending at regular intervals – usually on a ten-day cycle – to a
363 depth of 2000m and then ascending to the surface taking temperature and salinity measurements.
364 Since around 2007, Argo floats have provided quasi-global sampling of the oceans. The
365 temperature sensors are calibrated before the float is released and the manufacturer’s stated
366 accuracy is 0.002°C with a stability of 0.0002°C/year (Abraham et al. 2013). Whether this
367 accuracy is realized in the field is difficult to assess, but floats that have been recovered
368 remained within the manufacturer’s stated limits (ibid).

369

370 The good coverage of Argo floats, combined with the accuracy of the measurements they make
371 presents a dilemma for the data set creator. On the one hand, it would seem sensible to use high
372 quality measurements like these in the analysis (e.g. Huang et al. 2017). On the other hand, there
373 is much to be gained from using the Argo measurements as independent validation. The latter
374 approach is common in the satellite SST community (see e.g. Merchant et al. 2012, Berry et al.
375 2018), where drifting buoy data are often used for calibration and thus cannot be used for

376 validation. We adopt the same approach and reserve Argo for validation of the final product. The
377 relatively infrequent sampling provided by Argo – one profile every ten days – when compared
378 to say drifting buoys, which provide one measurement every hour, or even ships, which usually
379 measure once every six hours, means that including Argo in our analysis would have a relatively
380 small impact on the gridded anomalies. However, including Argo leads to a reduction of
381 measurement and sampling uncertainty of up to 30%.

382

383 **2.5.2 ATSR**

384 The ATSR (Along Track Scanning Radiometer) Reprocessing for Climate (ARC) data set
385 (Merchant et al. 2012) is a “climate-quality” analysis of SST retrievals from the ATSR satellite-
386 based instruments. The ATSR instruments were designed to make climate quality measurements
387 of SST and had a number of features to help achieve this. First, each instrument had a blackbody
388 onboard that allowed for continual calibration checks. Second, the satellites had a dual view
389 configuration, with observations made directly downwards (Nadir) and forwards (55° off
390 vertical) relative to the satellite’s motion. The two views allow the satellite to observe the same
391 area of the surface via two different paths through the atmosphere. By comparing the two, it is
392 easier to identify contamination arising from dust in the air or sulfurous volcanic particles. The
393 instruments also had three infrared channels, which allow more sensitive SST retrievals than is
394 possible using the two-channels available on the Advanced Very High Resolution Radiometer
395 (AVHRR) instruments (Merchant et al. 2014).

396

397 The ARC project reprocessed SST retrievals from the ATSR instruments. The SSTs we use here
398 are representative of a nominal measurement depth of 20cm (see Section 2.1) and an observation

399 time of 10:30 am and pm local time. Comparisons with Argo and drifting buoys show that there
400 is a minimal residual bias at a global scale (although locally there are deviations of order 0.1°C),
401 and that the uncertainty estimates provided with the data set are reliable. Three-way comparisons
402 with other satellites and buoys show that the individual ATSR retrievals have a typical
403 uncertainty of around 0.15°C (O'Carroll et al, 2008, Lean and Saunders 2013). The stability of
404 the ARC record has been demonstrated in the Tropical Pacific by comparison to moored buoys
405 in the Tropical Atmosphere-Ocean (TAO) array to be of order $0.01^{\circ}\text{C}/\text{decade}$ (Merchant et al.
406 2012).

407

408 The production of the ARC data set is almost entirely independent of SST measurements made
409 *in situ*. There is an indirect dependence as ARC uses reanalysis profiles to estimate the optimal
410 retrieval coefficients and the reanalysis used is driven using SST data sets that incorporate *in situ*
411 measurements.

412

413 We use the more reliable dual-view three-channel retrievals of SST from Version 1.1.1 of the
414 ARC data set, which span the entire period of the ATSR record 1991-2012. However, the period
415 from 1991 to 1995 is only intermittently covered due to the failure of one of the infrared
416 channels on the ATSR1 instrument shortly after launch.

417

418 **2.5.3 Buoys**

419 Although drifting and moored buoys are not independent of the data set that we develop in this
420 paper, they do consist of a single type of instrument or closely-related instruments and they are
421 of demonstrably higher quality than ship data (Kennedy 2014). There have been changes to

422 drifting buoy design over the years, but the largest changes had occurred by the early 1990s.
423 There is still some diversity in the design of drifting buoys and they are produced by a number of
424 different manufacturers. The nominal measurement depth varies, but is typically in the range 20-
425 50cm. Most sensors are of a nominal 0.1°C accuracy (Sybrandy et al. 2008), but estimates of
426 measurement uncertainty made in the field vary somewhat. Nevertheless, the buoy record is
427 considerably more homogeneous and stable than the unadjusted ship record over the period 1991
428 to present.

429

430 Moored buoys come in a variety of forms. The measurement depth is typically around 1m.
431 Sensor accuracy is also variable. Representative uncertainty values are given in Table 1. Some
432 moored buoys perform better than the average. In particular, measurements from the moored
433 buoys in the tropical Pacific from the TAO/TRITON array are of generally higher quality than
434 the US coastal arrays (K11c).

435

436 **3 Gridding and basic data preparation**

437 We averaged the SST measurements from the individual reports onto a 5° latitude by 5°
438 longitude monthly grid in a two-step process (R06). First, we sorted the observations into 1°
439 latitude by 1° longitude by pentad bins. Each SST observation was then converted into an
440 anomaly by subtracting the climatological average (for the period 1961-1990) for that 1° latitude
441 by 1° longitude pentad bin. We rejected anomalies with magnitude exceeding 8°C and calculated
442 the Winsorised (a form of trimmed mean, see Wilcox 2001) average of the remaining anomalies.
443 We then sorted the 1° latitude by 1° longitude pentad “super observations” into 5° latitude by 5°
444 longitude pseudo-month bins and took the Winsorised averages of the super-observations in each

445 of the larger bins. Figure 1(b, d, f) shows the number of “super observations” available per
446 month.

447

448 The contribution of each observation to the grid-box average has a weight, w , equal to

$$w = \frac{1}{ab} \quad \text{Equation 3.1}$$

449 where a is the number of observations in the same super observation and b is the number of
450 super observations in the larger 5° pseudo-month bin. The sum of the weights of all observations
451 in a 5° grid box equals one. Using these weights, we calculated the fractional contributions of
452 different measurement methods to the grid-box average and used these to estimate the fractional
453 contribution of each measurement method to the global and hemispheric averages (Figure 2).
454 Before 1915, bucket measurements have a weight of one in the global average. The weighting
455 does not perfectly reflect the influence of a single observation on the average because of the
456 Winsorisation process.

457

458 There are interesting changes in the influence of different measurement methods on the global
459 and hemispheric averages on all time scales. The Second World War stands out because of the
460 abrupt changes at the start and end, but there are rapid changes at other times such as the early
461 1960s (note that measurements of unknown type were assumed to be from ships). From around
462 1980 to 2005, the number and influence of buoys increases steadily after a brief peak in the late
463 1970s arising from a mass deployment of buoys during the First GARP (Global Atmospheric
464 Research Project) Global Experiment (FGGE, Garrett 1980). Despite the large numbers of
465 observations made by buoys since then, amounting to around 90% of all observations in the past
466 ten years (see e.g. Woodruff et al. 2011), the influence on the global average (i.e. the area

467 average of w for buoys as shown in Figure 2) in HadSST.4.0.0.0 has not consistently exceeded
468 50%. Although buoys make large numbers of observations, they typically do so in a limited area.
469 Observations from ships tend to be fewer, but more widely spread so the effective sampling-per-
470 observation is higher for a ship than it is for a buoy, particularly a moored buoy. This can also be
471 seen in the count of super-observations (Figure 1) which is largest in the 1970-1990 period when
472 the VOS fleet was at its peak.

473

474 In addition to the main gridded data set, we also gridded a number of subsets of the data which
475 we used to estimate biases for different measurement methods:

- 476 1. moored and drifting buoy observations;
- 477 2. all ship measurements;
- 478 3. measurements likely to have been ERI measurements (observations identified as an
479 engine room measurement, or taken by a ship from a country where more than 90% of
480 the fleet used engine room measurements at that time);
- 481 4. measurements likely to have been bucket measurements using a similar criterion;
- 482 5. hull sensor measurements; and
- 483 6. observations to which we could not attach a near-definitive measurement method.

484 **3.1 Uncertainty estimation for gridded averages**

485 Consider an SST measurement O_{ij} taken by agent i (either a ship or buoy) at space-time point j to
486 be the combination of the true SST, T_{ij} , together with a set of error components. The three error
487 components are: U_{ij} associated with uncorrelated errors, B_{ij} associated with “micro-bias” errors
488 that were correlated for all measurements made by the same agent, but uncorrelated between

489 agents, and G_{ij} a “macro-bias” error common to and correlated across all agents of a particular
 490 kind e.g. all ships taking ERI measurements.

$$O_{ij} = T_{ij} + U_{ij} + B_{ij} + G_{ij} \quad \text{Equation 3.2}$$

491 We calculate uncertainties in the gridded averages, σ_{gridbox} , using a variation of the formula from
 492 K11a,

$$\sigma_{\text{gridbox}}^2 = \sum_{i=1}^m \sum_{j=1}^{n_i} w_{ij}^2 \sigma_{u_i}^2 + \sum_{i=1}^m w_i^2 \sigma_{b_i}^2 + \sum_{h=1}^o w_h^2 \sigma_{G_h}^2 + \sigma_{\text{sampling}}^2 \quad \text{Equation 3.3}$$

493 where w_{ij} is the weight (as defined in Equation 3.1) in the grid-box average of the observation at
 494 point j made by agent i in that grid box. There are a total of m agents taking n_i measurements in
 495 the grid box. σ_{u_i} is the uncertainty associated with uncorrelated measurement errors and σ_{b_i} is
 496 the uncertainty associated with correlated “micro-bias” measurement errors for agent i . σ_{G_h} is
 497 the uncertainty associated with large-scale correlated errors from a particular measurement
 498 method, h , of which there are o different types and w_h is the weight assigned to measurements
 499 made using method h in that grid box. The values used for σ_u and σ_b are given in Table 1 and
 500 come from K11c. w_i is the sum

$$w_i = \sum_{j=1}^{n_i} w_{ij} \quad \text{and} \quad \sum_{i=1}^m w_i = 1 \quad \text{Equation 3.4}$$

501

502 The sampling uncertainty, σ_{sampling} is given by:

$$\sigma_{\text{sampling}}^2 = \frac{1}{n_s} \sigma_s^2 [1 - \bar{r}] \quad \text{Equation 3.5}$$

503 Where n_s is the number of super observations in the 5° pseudo-month grid box. σ_s is the variance
 504 at a space-time point within the grid cell (here assumed not to vary across the grid cell) and \bar{r} is
 505 the average correlation of space-time points within the grid cell (for details of the calculation see

506 K11a). We use the number of super observations rather than the number of observations as some
 507 agents such as moored buoys make many hundreds of observations in a single location and it is
 508 clear that the sampling uncertainty cannot be endlessly reduced by making more observations in
 509 the same place. A count of super observations gives a better idea of the number of independent
 510 space-time points sampled, though it remains less than perfect. The sampling uncertainty
 511 parameters were estimated as in K11a but using the 5° gridded ICOADS release 3.0 data
 512 between 1961 and 2016.

513 3.2 Correlated errors

514 The covariance, $R(p, q)$, associated with correlated “micro bias” errors from the same ships
 515 visiting grid boxes p and q was

$$R(\mathbf{p}, \mathbf{q}) = \sum_i \mathbf{w}(\mathbf{p})_i \mathbf{w}(\mathbf{q})_i \sigma_{b_i}^2 \quad \text{Equation 3.6}$$

516 where i sums over all the ships that visited grid boxes p and q .

517

518 Correlated errors between grid-box averages can also occur when two different ships using the
 519 same measurement method visit different grid cells. In estimating the biases and uncertainties for
 520 ERI data (see Appendices A1 and A2) we assumed that G_{ij} (Equation 3.2) was non-zero and
 521 equal for all ships making ERI measurements. The covariance is calculated in a similar way to
 522 that of the micro-bias errors above.

$$R(\mathbf{p}, \mathbf{q}) = \sum_h \mathbf{w}(\mathbf{p})_h \mathbf{w}(\mathbf{q})_h \sigma_{G_h}(\mathbf{p}) \sigma_{G_h}(\mathbf{q}) \quad \text{Equation 3.7}$$

523 where $\sigma_{G_h}(\mathbf{p})$ is an estimate of the uncertainty in the large-scale bias associated with

524 measurement method h at location p . For ERI measurements, $\sigma_{G_h}(\mathbf{p})$, was assumed to be

525 constant, but for bucket measurements, $\sigma_{G_h}(p)$ was assumed to be equal to the bucket adjustment
526 field at that location. In other words, we assume that the shape of the bucket adjustments is
527 known but not the exact magnitude. In effect, the biases are estimated by regression (Appendices
528 A1 and A2), with the covariances defined in Equation 3.6 and 3.7 specifying both the shape of
529 the predictor and its prior variance.

530 **3.3 Use of the error model**

531 We use the error model in two distinct ways in this paper. First, it is used to estimate some of the
532 error terms described in Equation 3.2, that is to determine the actual sizes of biases such as G_{ij} .
533 The details of this are given in Appendix A1 and A2. Second, the error model is used to estimate
534 uncertainties in derived quantities such as grid box averages (Equation 3.3) as well as regional
535 and global averages. In order to calculate regional and global averages, Equation 3.3 and
536 Equation 3.6 are combined to produce a total error covariance matrix and the uncertainties are
537 propagated using the propagation of uncertainties formula for a weighted average where the
538 weights are the grid-cell areas. The specific form is described in Section 3.3 of K11a.

539 **3.4 Example fields**

540 Figure 3 shows example fields for June 2003. The number of observations exceeds 100/month in
541 many grid boxes and is somewhat homogeneous thanks to the widespread use of drifting buoys.
542 In contrast, the number of super observations varies greatly, with the highest numbers in the
543 northern hemisphere shipping lanes, demonstrating again the more efficient sampling per
544 observation afforded by ships. The very localized sampling of moored buoys is also clear.

545

546 The uncorrelated error component of the uncertainty is relatively small except for at the edges of
 547 the observational coverage, where there are very small numbers of observations contributing to
 548 each grid-box average. Uncertainties associated with micro-bias errors are typically larger as
 549 they are related to the number of platforms contributing to the average. They are lowest over the
 550 north Pacific and North Atlantic where there is a great diversity of measurement platforms.

551

552 The greater weight given to ship observations, which are generally less reliable, is potentially a
 553 weakness in the simple approach adopted here. A more sophisticated method such as an optimal
 554 interpolation scheme (Karspeck et al. 2012) might give a greater weight to the more reliable
 555 buoy observations. However, the current method uses Winsorisation, which minimize the effect
 556 of outliers and makes it much easier to keep track of the correlation structure of the errors
 557 through the full uncertainty calculation.

558

559 **4 Bias Adjustments**

560 The estimation of biases has two basic steps. First, we create grids of the fractional contribution
 561 of each measurement type to the gridded averages. Second, we assign biases to each
 562 measurement method and calculate an overall bias in each grid cell. The bias B in a grid cell is
 563 equal to

$$B = f_e E + f_c B_{tc} + f_w B_{tw} + f_r B_{tr} + f_d D \quad \text{Equation 4.1}$$

564 Where f_e is the fractional contribution of ERI and hull sensor measurements, f_c is the fractional
 565 contribution of measurements made with canvas, or otherwise-uninsulated, buckets; f_w is the
 566 fractional contribution from wooden buckets; f_r is the fractional contribution of measurements
 567 made with rubber, or otherwise-insulated modern buckets; and f_d is the fractional contribution of

568 measurement from buoys. E , B_{tc} , B_{tw} , B_{tr} and D are the biases associated with each of these
 569 measurement methods respectively.

570

571 In practice, we do not have estimates of the true historical bucket biases, B_{tc} and B_{tw} (the t in the
 572 subscript stands for “true”). What we do have are the corrections B_c and B_w (from FP95) which
 573 adjust a bucket measurement to be consistent with the average bias in the period 1961 to 1990
 574 (rather than adjusting a bucket measurement to the true SST). To convert B_c to B_{tc} , it is necessary
 575 to calculate:

$$B_{tc} = \frac{B_c + \overline{f_r} B_{tr} + \overline{f_e} E}{(1 - \overline{f_c})} \quad \text{Equation 4.2}$$

576 where an overbar denotes the average for that value over the period 1961-1990. The derivation,
 577 which is a modified from of that in K11b, can be found in Appendix A4. Unlike in K11b, we will
 578 make a direct estimate of B_{tr} from the data (Section 4.1.3) so the formula shown here is
 579 somewhat different.

580

581 Another difference from K11b is that we assume D is zero. In other words, we believe that buoys
 582 provide both an accurate measurement of SST and a benchmark (as noted in Section 2.1) for
 583 estimating the biases associated with other measurement methods. A similar approach, using
 584 drifting buoys as a baseline, is taken in ERSSTv5 (Huang et al. 2017) and drifting buoys are
 585 widely used as “ground truth” in the satellite SST community (e.g. Lean and Saunders 2013 and
 586 Embury et al. 2012). In contrast, in K11b it was assumed that the time series E was known and
 587 the adjusted ERI measurements were used as the baseline. A similar approach was taken for
 588 ERSSTv4. Note that the choice of whether to use ships or drifters as a baseline, in so far as this
 589 constitutes a constant offset, does not affect the estimation of trends or anomalies. It can,

590 however, affect the estimation of actual SSTs, which are important for some applications. It also
 591 contributes to a more intuitive presentation of the time series, with smaller uncertainties in the
 592 drifter-rich period (Figure 8).

593

594 In practice, when we come to estimate the biases in the data where we do not know the
 595 assignments perfectly, the bias in an average of a group of observations will be some linear
 596 combination of biases from both bucket and ERI measurements. For example, if we compare a
 597 collection of observations labeled as buckets to a set of unbiased drifting buoy data then the
 598 empirical bias B^* seen in the bucket-labelled data (in this case assuming that all the buckets are
 599 rubber) would be:

$$B^* = f_{correct}B_{tr} + (1 - f_{correct})E \quad \text{Equation 4.3}$$

600 where $f_{correct}$ is the fraction of correctly-labeled bucket observations.

601 **4.1 Implementation**

602 In order to calculate biases in the data using Equation 4.1 described above, we need to assign
 603 values to each of the components. In the following subsections, we describe how the biases are
 604 estimated and how we refine some of the metadata assignments. Due to uncertainty in many of
 605 those values and the complicated interactions and correlations between them, we take an
 606 ensemble approach and generate many different sets of possible biases that span some part of the
 607 overall uncertainty in the adjustments. The process of estimating the biases is broken down into a
 608 number of steps:

- 609 1. Generate a set of correction fields for canvas and wooden buckets (B_c and B_w) using a
 610 modified version of R06 (see Section 4.1.1) and SR02. These will form the basis for the
 611 bucket adjustments in the early part of the record.

- 612 2. Generate an ensemble of estimates of $E(t)$, the time-and-space-varying biases associated
613 with ERI measurements, and B^* (Equation 4.3), the bias associated with observations
614 labeled as buckets (see Section 4.1.2). Note that the observations labeled as buckets will
615 contain some unknown fraction of mislabeled ERI measurements. These are dealt with in
616 step 4.
- 617 3. Estimate the spatially- and seasonally-varying biases for modern insulated buckets and
618 generate an ensemble of bucket biases (B_{tr}) (see Section 4.1.3).
- 619 4. Using the values generated in steps 1-3, generate an ensemble of start and end dates for
620 the transition from canvas to insulated buckets and estimates of the fraction of
621 observations that are correctly identified as bucket measurements, $f_{correct}$ (see Section
622 4.1.4). The ensemble of estimates of $f_{correct}$ gives a measure of the uncertainty in the
623 metadata assignments.
- 624 5. Generate an ensemble of estimates of how unknown measurement types are to be
625 assigned (Section 4.1.5) and a separate ensemble of parameters for the Second World
626 War, which reflects the fact that shipping and behaviors changed during the war and are
627 somewhat independent from and more uncertain than the periods before and after (see
628 Section 4.1.6).

629 We describe these steps in the following subsections and a schematic representation is shown in
630 Figure 4. Comparisons with sub-surface profile data were used to inform or constrain some of
631 the parameter ranges. In order to minimize the effect of over-fitting, particularly when data were
632 few, the constrained parameters were not tied too closely to each other and were chosen to
633 represent a reasonable range of uncertainty.

634 **4.1.1 Generating wooden and canvas bucket correction fields, 1850-1941**

635 The method described in K11b (in turn a modified version of R06) was the basis for generating
636 an ensemble of 200 bucket correction fields for the data prior to 1942. A number of changes
637 were made to the processing where uncertainties were likely to have been underestimated in the
638 previous version.

639

640 First, we treated errors in the monthly correction fields as correlated: a single number was drawn
641 from a standard normal distribution, multiplied by the uncertainty from R06 and applied to all
642 monthly fields for a particular ensemble member. In K11b, a separate draw was made for each of
643 the twelve months, effectively treating the errors as uncorrelated. However, in the original
644 version of the method (R06) although they were independent draws and hence uncorrelated, the
645 95% uncertainty ranges were calculated separately for each month and then combined. This step
646 effectively treated them as if the errors were fully correlated. This change brings the two
647 methods back into line and increases the uncertainty in the bucket biases at annual and longer
648 time scales relative to K11b.

649

650 Second, in the calculation of the difference in annual tropical average SST and NMAT – used to
651 fix the fractions of wooden and canvas buckets from 1850 to 1920 – the estimated uncertainties
652 arising from measurement errors were increased to account for the correlated micro-bias errors.
653 We calculated the uncertainties assuming the errors were uncorrelated and then multiplied the
654 resulting uncertainties by 5.08 (using the conversion factor for the tropics from K11a which
655 accounts for the spatial correlation of the errors ($= 2.2$) multiplied by the square root of $12/2.25$
656 for the temporal correlations). This again has the effect of increasing the uncertainty in the

657 estimated bucket biases prior to 1920 relative to K11b and R06 where these errors were treated
658 as if they were uncorrelated.

659

660 Third, we generated half of the 200-member ensemble as in R06, by assuming a linear transition
661 from wooden to canvas buckets. The other 100 ensemble members were generated by assuming
662 a step-change with different constant fractions of wooden and canvas buckets before 1906 and
663 after 1910, with a linear change between these states from 1906 to 1910. The step change
664 provides a qualitatively better fit to the noisy tropical SST-NMAT temperature differences used
665 to estimate the transition (not shown). Rather than being a consequence of a step change in the
666 fractions of canvas and wooden buckets, the step change in the SST-NMAT difference might
667 instead be due to an increase in ship speeds around this time (Carella et al. 2017a). The change
668 that we apply would be about the same in either case because an increase in the speed of the ship
669 would also increase the necessary correction.

670

671 Fourth, every even ensemble member (2, 4 ... 198, 200) was a blended average of R06 style
672 adjustments and SR02 style adjustments (Equation 4.4). SR02 use patterns of SST-NMAT
673 differences to estimate their adjustments. For our SR02-style adjustments, we estimated fields of
674 adjustments associated with canvas buckets by taking the SR02 corrections for 1941. Wooden
675 bucket corrections were calculated by assuming that the 1850 SR02 corrections were 80%
676 wooden buckets and 20% canvas buckets. We estimated an uncertainty of 10% for the SR02-
677 style bucket corrections i.e. the bucket corrections were multiplied by a number (A_0) drawn from
678 a normal distribution with mean 1 and standard deviation of 0.1. For every even ensemble
679 member, a weighted average of the R06 and SR02 adjustments was taken with the weight of the

680 R06-style correction, A_1 drawn from a uniform distribution in the interval [0,1] and the weight of
 681 the SR02-style correction equal to $1-A_1$. In the extreme cases, the corrections look either entirely
 682 like R06, or entirely like SR02. For odd ensemble members, A_1 was set to 1. For each ensemble
 683 member a new value for B_{cR06} was drawn and combined like so:

$$\mathbf{B}_c = A_1 \mathbf{B}_{cR06} + (1 - A_1) A_0 \mathbf{B}_{cSR02} \quad \text{Equation 4.4}$$

684 The use of SR02 adjustments in addition to R06-style adjustments means it is possible to use a
 685 parametric framework to explore some of the structural uncertainties (Thorne et al. 2005) that
 686 would ordinarily only be accessible by comparing different data sets. The SR02 adjustments tend
 687 to follow the climatological pattern of sensible heat fluxes, whereas the R06 adjustments more
 688 closely follow latent heat fluxes. By using a weighted average of the two for half of the ensemble
 689 members, we explore different combinations of the two very different approaches and hence a
 690 wider spectrum of adjustments with different relationships to specific and latent heat fluxes.
 691 Time series of the resulting corrections are shown in Figure 5.

692

693 **4.1.2 Estimating biases for individual measurement methods, 1940 onwards**

694 We generated a subset of the data after 1930 using only measurements that we identified as ERI
 695 or that had a fractional assignment to ERI greater than or equal to 0.9. Error covariances were
 696 estimated as described in Section 3, Equation 3.6. The uncertainty associated with the large-scale
 697 correlated errors arising from using ERI measurements, σ_G , was set to be 0.2°C (the estimated
 698 mean bias from K11b). Monthly ERI biases were then estimated by comparison with drifting
 699 buoys and six different versions of the near-surface profile data for 1930-2014, from
 700 HadIOD.1.2.0.0, using the method described in the Appendix (Section A2). Figure 6 shows the

701 global average of the combined large-scale systematic errors and micro-bias errors estimated in
702 this way.

703

704 The left hand column of Figure 6 shows ERI biases estimated from comparisons with the six
705 different data sets created from near-surface profile data for the globe, southern hemisphere and
706 northern hemisphere.

707

708 ERI biases are generally positive apart from in the very earliest years when uncertainties are
709 large. The largest reliably-estimated ERI biases occur between 1955 and 1970 when they range
710 from 0.2 to 0.6°C. After 1970, ERI biases drop to a local minimum in the mid 1990s. They then
711 rise again to a peak in the early 2000s before dropping again, approaching 0°C around 2018. The
712 reasons for the variations in ERI bias with time are not clear, but the fleet of ships making ERI
713 measurements was not designed to make climate quality measurements so the instability itself is
714 unsurprising.

715

716 Measurements labeled as buckets and hull sensors were analysed in the same way. For buckets
717 $\sigma_G(p)$ was set to the value of the bucket correction field at location p . The bias estimation is
718 effectively a Bayesian regression using the bucket correction field as a predictor. For hull
719 sensors, σ_G was set to 0.2°C as it was for ERI measurements. Bucket biases (see left-hand
720 column of Figure 6) are rather variable but less so than in the ERI set. The bias is often positive,
721 which is contrary to the general expectation that, on average, buckets lose heat (Kent et al. 2017).
722 However, buckets can exhibit a warm bias during the day if the sun shines on the bucket and
723 warms the water sample, and that might offset heat losses at other times (although in FP95, the
724 solar effects in the correction fields were rather small). It is also likely that there is some

725 contamination of the bucket subset by mislabeled ERI measurements. These issues are dealt with
726 further in Section 4.1.3 and 4.1.4. Hull sensors show biases very similar to those of ERI
727 measurements from 1990 onwards, so we combined these two data sources in the bias
728 adjustments.

729

730 The ERI biases in each grid box were smoothed in time, but not in space. We generated an
731 ensemble of 200 sets of adjustments. For each ensemble member, we randomly selected one of
732 the six bias-adjusted sub-surface data sets to use as a basis, drew samples from the posterior
733 covariance of the estimated ERI biases and added these to the mean biases (see section A2). The
734 samples for each grid box were then smoothed in time using a LOWESS (LOcally WEighted
735 Scatterplot Smoother) filter (<http://flux.aos.wisc.edu/data/code/idl-lib/util/bueilib/lowess.pro>)
736 with a width of 121 months (± 5 years). Because the ERI biases can only be estimated where
737 there are known ERI measurements, there can be long temporal gaps in individual grid-box
738 series. Where fewer than 30 months of data were available for a particular grid box in a 121-
739 month period centred on a particular month, the grid box value for that month was set to the
740 global average ERI bias. This procedure preserves the spatial structure seen in the ERI biases and
741 their spatial covariances at the expense of some loss in temporal resolution. Spatial structure
742 could arise for a number of reasons; for example, the biases of ships travelling along a common
743 shipping lane could be quite different to the biases of ships travelling in the same region but not
744 following the same routes.

745 **4.1.3 Estimating parameters for modern bucket biases, 1970s onwards**

746 Carella et al. (2018) used ICOADS flags SI (SST measurement method) and SIM (SST
747 measurement method from WMO Pub. 47) to identify bucket measurements. They found that

748 their compound flag, SI(M), was confirmed in 90% of cases from 1970 onwards. We selected
749 bucket observations using SI, or SIM where SI was missing or ambiguous, and gridded these
750 separately. This is similar to the assignment made in the previous section, but does not include
751 fractional assignments. Biases were estimated for this data set relative to the sub-surface data
752 from 1970-2005 using the method described in the Appendix (Section A2). For each calendar
753 month, we computed fields of the mean bias and its standard deviation and the fields were filled
754 using simple Kriging (Clark 1979) with a fixed angular length scale of 15° in latitude and
755 longitude. 200 different versions of the adjustments were calculated by drawing samples from
756 the posterior distribution of the Kriged solution.

757

758 The resulting fields show a small overall warm bias (Figure 5d) and there is a strong seasonal
759 cycle in the northern hemisphere, which peaks in the summer months (not shown). This suggests
760 that there is a solar heating bias in modern bucket measurements that exceeds the small heat
761 losses found at night under moderate wind conditions by Kent and Kaplan (2006). The warm
762 bias primarily affects the high latitude oceans (poleward of 40°N and S) and there is a small cold
763 bias throughout the tropics. This pattern is similar to that seen by Carella et al. (2018) in which
764 bucket measurements were biased warm relative to ERI measurements at high latitudes from the
765 late 1980s onwards. The pattern identified by Carella et al. (2018) peaked towards local noon,
766 but it is not clear whether the warm bias is due to the bucket sampling a shallow surface layer, or
767 due to direct solar heating of the bucket itself. In either case it is distinct from buoy
768 measurements which are the target for our analysis (Section 2.1)

769 **4.1.4 Generating start and end dates for the transition from canvas to modern buckets and**
 770 **estimating the fraction of measurements misidentified as buckets**

771 In the middle to late twentieth century there was a transition from the use of canvas buckets to
 772 the use of rubber buckets. This transition occurred at the same time as large numbers of ships
 773 adopted the ERI method. Neither transition is well documented and metadata, particularly that
 774 implied by deck and country information (Section 2.2.1), is not completely reliable. We can use
 775 the estimated biases from the previous sections (Section 4.1.1 to 4.1.3) together with proposed
 776 start and end dates for the canvas to rubber bucket transition to infer the fraction of incorrect
 777 metadata. Where this fraction takes an impossible value – outside the range [0,1] – we can reject
 778 that combination of start and end dates and thus narrow the uncertainty range for these
 779 parameters. Once we have narrowed down the range of start and end dates in this way, we can
 780 generate a best estimate of the time series of the fraction of correct metadata and a plausible
 781 range within which it can be varied to generate an ensemble.

782

783 We started with a wide range of dates for the transition from canvas to insulated buckets. Start
 784 dates were initially in the interval 1930 to 1960 and end dates between one year after the start
 785 date and 1980. The transition from canvas to insulated buckets was assumed to be linear between
 786 the start and end dates and the same everywhere.

787

788 Similarly to Equation 4.3, if we assume that B^* , the bias estimated from observations labelled
 789 (perhaps incorrectly) as buckets, is a combination of biases associated with buckets, B , and
 790 biases associated with ERI, E , then we can write B^* as

$$B^* = f_{correct}B + (1 - f_{correct})E \quad \text{Equation 4.5}$$

$$B = f_{canvas}B_{tc} + (1 - f_{canvas})B_{tr} \quad \text{Equation 4.6}$$

791 where $f_{correct}$ is the unknown fraction of measurements labeled as bucket measurements that were
 792 correctly identified, f_{canvas} is the fraction of canvas buckets changing linearly from one to zero
 793 between the start and end date of the transition. Rearranging, we get

$$f_{correct} = \frac{B^* - E}{B - E} \quad \text{Equation 4.7}$$

794 With an estimate of B , B^* and E it is therefore possible to get an estimate of $f_{correct}$. B_{tc} was
 795 estimated by subtracting the climatological-average ship bias from the mean bucket correction B_c
 796 (Section 4.1.1) with the ship bias being calculated relative to each of the six different sub-surface
 797 data sets (Section 2.3) and then averaged across them. B_{tr} was calculated as in Section 4.1.3.

798

799 For each pair of start and end dates the estimates of B^* , E and B , for each month were used to
 800 derive an estimate of $f_{correct}$ for each month. Note that values of $f_{correct}$ slightly in excess of one
 801 are possible due to measurement errors in B^* , E or B . A simple uncertainty range on $f_{correct}$ was
 802 estimated by increasing and decreasing the bias of rubber buckets by 0.05°C (a lower bound on
 803 the uncertainty of the method, see Appendix A3.3). If the uncertainty range in annual average
 804 $f_{correct}$ did not overlap the range $[0.5, 1.0]$ during the transition period from 1955 to 1962, the start
 805 and end dates were rejected. In practice, no combination was consistent with a value less than
 806 0.5. The period 1955 to 1962 was found to be particularly sensitive to the choice of start and end
 807 points (the spread in Figure 7a and b is particularly wide during this period) while also having
 808 reliable sub-surface data for estimating the biases. Prior to 1955, data coverage of the sub-surface
 809 data is much more sparse and therefore the uncertainties are larger and provide a much less
 810 useful constraint. After 1962, the spread is already well-constrained and nothing is gained by
 811 extending the constraint period beyond this point.

812

813 Figure 7 shows the accepted and rejected start and end dates for the transition and the associated
814 time series of $f_{correct}$. Few start dates were rejected by this method, but end dates for the transition
815 before 1961 were inconsistent with the data. This still gives a wide range of possible start and
816 end dates, including those which describe a relatively rapid transition starting in the late 1950s.

817

818 The time series of $f_{correct}$ is broadly consistent with an independent assessment of the metadata
819 made in Carella et al. (2018). Metadata are largely reliable from the early 1980s to the early
820 1990s. Another period of higher reliability is seen from around 1955 to 1962. Between 1965 and
821 1980, there is a period of less reliable metadata. Carella et al. (2018) argue for a relatively rapid
822 transition from uninsulated to insulated buckets between the mid 1950s and mid 1960s, which is
823 consistent with the results shown (Figure 7c). For example, such a transition starting in 1955 and
824 ending in 1965 is accepted.

825

826 As with the ERI adjustments, a set of perturbed realisations of $f_{correct}$ were calculated that were
827 similar to, but not tightly constrained by the best estimate. We calculated a time series of central
828 values for $f_{correct}$ by averaging the time series of $f_{correct}$ across all allowed pairs of start and end
829 dates (Figure 7d). Before 1952, the best estimate for $f_{correct}$ was set to 0.5 and after 1978 it was
830 set to 0.95. This gives a continuous series that was then smoothed with a LOWESS filter with a
831 width of 4 years. The one-sigma uncertainty was set to: 0.15 before 1952; 0.1 from 1952 to 1978;
832 and 0.05 from 1978 onwards. This uncertainty range encompasses the majority of annual
833 averages except in the post 1978 period. After 1978, there are some sharp variations, but these

834 occur during a period when bucket measurements made a smaller and smaller contribution to the
835 global average.

836

837 200 random time series with a lag-1 correlation of 0.99 and values in the range [-1,1] were
838 generated, scaled by the estimated uncertainty in $f_{correct}$ and added to the best estimate. Values
839 were capped at one with generated values above one set to one. Several example series are
840 shown in Figure 7d.

841

842 **4.1.5 Measurements with unknown method**

843 Some measurements cannot be assigned a measurement method (Figure 2). A fraction of these
844 unknown measurements was randomly reassigned to be either bucket or engine room
845 measurements. A monthly time series was created which varied randomly between zero and one
846 with an autocorrelation of 0.99 as in K11b. For a given month, the contribution of unknown
847 measurements in each grid box was multiplied by this number and added to the contribution
848 from ERI measurements. The remainder was added to the contribution from bucket
849 measurements.

850

851 **4.1.6 Parameters for the Second World War, 1941-1945**

852 During the Second World War, there was widespread disruption to shipping. There are also
853 discontinuities in ICOADS data sources and rapid changes in bias at the start of the war (R06,
854 Huang et al. 2017) and at the end (Thompson et al. 2008). FP95 suggested that the Second World
855 War saw the hasty (and even permanent) adoption of ERI measurements, as they were safer to

856 make. The rapid changes seen in Figure 2 during this period are also partly due to our choices for
857 the initial metadata assignment. In order to reflect the greater uncertainty during the war years,
858 we generated a separate set of parameters with increased uncertainty for the period between
859 January 1941 and August 1945.

860

861 The fraction of measurements that were labeled as buckets, but which were really ERI, was set to
862 a value selected from a uniform distribution in the range [0, 1] to reflect the possibility that ERI
863 measurements were temporarily taken in preference to bucket measurements across a large
864 fraction of the fleet (Figure 7d). The ERI bias was chosen from a uniform distribution in the
865 range [0, 0.5] and the fraction of unknown measurements set to buckets was drawn from a
866 uniform distribution in the range [0, 0.25]. These values were chosen to give a broad range of
867 possibilities – and hence large uncertainty in the adjustment – and reflect the likely prevalence of
868 ERI measurements.

869 **5 Results**

870 Figure 8 shows the estimated biases in global and hemispheric averages for all data sources
871 including ships and buoys. The biases shown in (Figure 8 a,c,and e) are relative to our reference
872 buoy SST at a nominal depth of 20cm (Section 2.1). The bias relative to the average bias in the
873 climatology period, 1961-1990, is also shown (Figure 8 b,d,and f). The difference between the
874 two can be sizeable because of the large absolute biases in the climatology period when the vast
875 majority of measurements were made by ships, which were affected by warm biases associated
876 with ERI measurements.

877

878 From 1850 to 1939, the bias becomes increasingly negative (Figure 8 a,c,and e) reflecting the
879 transition from wooden to canvas buckets. This is accompanied by an increase in the seasonal
880 cycle of the biases, which is particularly clear in the Northern Hemisphere and arises from the
881 seasonal drivers of the bucket biases: sensible, latent and solar heat fluxes (FP95). Between 1935
882 and 1939, the bias becomes slightly less negative as ERI measurements start to appear in
883 ICOADS. Between 1939 and 1941, the bias becomes rapidly more positive as ERI
884 measurements, principally from US sources, enter the record in large numbers (see Figure 6 of
885 R06). From 1941 to 1945, during the Second World War, uncertainties are larger and there is a
886 net positive bias that reflects the assumed increase in the prevalence of ERI measurements. The
887 bias falls briefly in the post war years, reflecting a partial reversion to bucket measurements.
888 However, the change is not so marked as it was in HadSST.3.1.1.0 because the ERI biases are
889 here estimated to be larger than was previously assumed and the reliability of the initial method
890 assignments (Section 2.2.1) is estimated, at times, to be worse than was assumed in
891 HadSST.3.1.1.0. The bias increases to around 0.2 to 0.3°C between 1955 and 1970. During this
892 period, ERI biases remain high and, towards the end, many observations that were initially
893 flagged as buckets were reassigned to be ERI measurements. In HadSST.3.1.1.0, a constant
894 fraction ($30\pm 10\%$) of measurements initially labelled as buckets were reassigned to be ERI with
895 a different fraction in each ensemble member. In HadSST.4.0.0.0, as well as being different for
896 each ensemble member, the reassignment is time-varying and is constrained based on the biases
897 estimated from comparisons with the sub-surface oceanographic profiles (Figure 7, Section
898 4.1.4).
899

900 From 1970 to 1980, the average bias drops reflecting a reduction in ERI bias (see also Figure 6).
901 From around 1990, buoys start to have a significant effect on the record. This reduces both the
902 bias in the combined data set and the uncertainty in the bias. Despite the increasing prevalence of
903 buoys, changes in bias from 1990 to the present continue to reflect changes in the bias of the ERI
904 and hull sensor measurements, which are the predominant means by which SST measurements
905 were made from ships in this period and which still have a weight of around 50% in the global
906 average (Figure 2).

907

908 The fall in bias in the ship data since the early 2000s is not confined to ERI measurements and is
909 consistent with the change in bias estimated by Huang et al. (2015, 2017) using comparisons
910 with NMAT and buoy data. It is seen in VOSclim ships (<http://sot.jcommops.org/vos/vos.html>),
911 which are a select subset of the full VOS fleet for which better quality metadata are available. It
912 is also seen in measurements from ships that are not part of VOSclim, and in both bucket
913 measurements and hull sensor measurements, which suggests that this is a pervasive reduction in
914 the overall bias in the ship data from the early 2000s to present.

915

916 The bias relative to 1961-1990 has a similar evolution (Figure 8 b,d,and f). However, there are
917 two things to note. First, as the absolute bias shows, there is a large warm bias in the data from
918 1961-1990 and the uncertainty is larger than it is for the later data. When biases in the buoy-rich
919 period from 1990 onwards are expressed relative to the 1961-1990 average, their uncertainty
920 increases. This happens because we are comparing an accurate measurement to an inaccurate
921 baseline and the overall uncertainty in the difference is therefore high. Second, the continued use
922 of buckets in the climatology period means that there is a residual seasonal cycle in the absolute

923 biases. When those biases are expressed relative to 1961-1990, the seasonal cycle is reduced
924 during the climatology period and in the pre-Second World War period, but increased outside of
925 it.

926

927 Figure 9 shows series of global and regional average SST anomalies for HadSST.4.0.0.0 and for
928 the unadjusted gridded measurements. The adjustments have the largest effect prior to the
929 Second World War. However, even after the war, the biases affect the long-term changes in the
930 data. The apparent step change in the mid-1970s is slightly larger in the adjusted data, with the
931 adjustments cooling the series in the late 1960s and warming it in the late 1970s. The difference
932 between the 2010s and the 1961-1990 period is also increased. This is caused by the overall
933 decrease in the bias over the past 50 years, which is principally due to the decreasing ERI bias
934 and the increasing influence of unbiased drifting buoys. There is a strong annual cycle in the
935 northern hemisphere average from around 2005 onwards. The origin of this is not clear, but may
936 be due in part to reduced sea-ice extent during summer months.

937

938 In the following subsections, we compare the new HadSST.4.0.0.0 data set and bias adjustments
939 with other marine temperature (SST and MAT) data sets to identify points of similarity and
940 difference. In Section 5.1 we compare subsets of data made using buckets and engine room
941 measurements and show that the adjustments reduce the relative bias between them. Section 5.2
942 details comparisons with instrumentally homogeneous data records from 1991 to 2018,
943 demonstrating the stability of the HadSST.4.0.0.0 record in the modern era. Comparisons with
944 longer MAT and near-surface oceanographic profile measurements series are described in
945 Section 5.3 and highlight an interesting discrepancy with HadNMAT.2.0.1.0. In Section 5.4 we

946 compare HadSST.4.0.0.0 to HadSST.3.1.1.0 and in Section 5.5 we extend the comparison to
947 other SST data sets showing good long term agreement between them.

948 **5.1 Internal consistency**

949 Two data sets were created. One contained only observations identified as likely bucket
950 measurements. The other contained only observations identified as likely ERI measurements.
951 The two data sets were gridded separately and their measurement and sampling uncertainties
952 estimated. Each of the two subsets of the data was bias adjusted using the relevant estimated
953 biases. Figure 10 shows the global average anomaly for the two collocated data sets before and
954 after adjustment.

955

956 Prior to adjustment, there is a clear time-varying offset between the two data sets between 1940,
957 the first year in which ERI measurements are present in significant numbers, and 2014. This is
958 consistent with various estimates (see K11b) of the relative biases between the two measurement
959 methods. There is similar inter-annual variability in the two data sets. The drop in the combined,
960 unadjusted SST anomalies in 1945 (Figure 8 and Figure 9) is not evident in either of the two
961 subsets. This suggests that the abrupt drop in estimated global temperature highlighted by
962 Thompson et al. (2008) is largely a result of a rapid change in the relative proportions of ERI and
963 bucket measurements in ICOADS at that point. On the other hand, the two subsets are each
964 noisier than the combined series and the discontinuity occurs at the point at which the estimated
965 uncertainty changes markedly so further analysis is warranted.

966

967 The adjustments improve the agreement between the two data sets. This indicates that the bias
968 adjustments are functioning as expected at a hemispheric scale. The agreement is expected

969 because adjustments to both strands of data and the parameter choices are set to ensure loose
970 consistency with the same reference data sets. The divergence post-2005 is likely due to the
971 scarcity of bucket data after this point, which is also reflected in the much larger estimated
972 uncertainties. The uncertainties in the adjusted data are also broader than in the unadjusted data.
973 This may seem counterintuitive and arises because bias errors are not included in the uncertainty
974 range for the unadjusted data. If they were, they might amount to a few tenths of a degree
975 reflecting the large differences between ERI and bucket data.

976 **5.2 Comparison with instrumentally homogeneous data sets**

977 In the modern period, the data set can be compared to three instrumentally homogeneous data
978 sets (Argo, ARC and buoys, see Section 2.2) that combine higher accuracy with good global
979 coverage. Hausfather et al. (2017) used similar instrumentally homogeneous data sets to assess
980 the effect of bias adjustments on trends in the modern period and detect drifts in the global
981 averages of combined series such as HadSST3 and ERSSTv3/v4 that were smaller than
982 0.05K/decade. Our comparison builds upon the Hausfather et al. (2017) analysis in two ways.
983 First, all the data sets are gridded using a common procedure (Section 3), which minimizes the
984 possibility of differences due to processing. Second, the same 1961-1990 climatology is used to
985 calculate anomalies for all of the data sets. Using a common climatology, rather than force each
986 series to average to zero across the recent period of overlap, means that any absolute differences
987 in SST will also be highlighted along with any trend differences.

988

989 We constructed gridded data sets of the three comparison data sources using the method
990 described in Section 3. Figure 11 shows global averages of the three data sets compared to

991 HadSST.4.0.0.0 at the locations where there is common data coverage. The agreement between
992 these data sets is very good with differences largely falling within the uncertainty ranges.

993

994 The least interesting comparison is between HadSST.4.0.0.0 and the buoys (Section 2.5.3) as
995 they are used in the construction of the HadSST.4.0.0.0 data set. Nevertheless, it is useful to
996 compare the two to ensure that we have not introduced a bias by the addition of the ship data. It
997 is clear that the adjustments greatly reduce the difference between the combined data set and the
998 buoys alone. There is, however, a period in 1992/1993 where the two diverge by much more than
999 the estimated uncertainty. The early 1990s are rather sparsely observed by drifting buoys and, in
1000 this case, the buoy record is likely to be erroneous with local large deviations associated with
1001 poor-quality drifting buoy data. Spatially (Figure 12(a)), we generally see small average
1002 differences (within $\pm 0.1^\circ\text{C}$ for most areas) except in coastal regions around the US and in the
1003 northern hemisphere western boundary currents.

1004

1005 A more stringent test is to compare the HadSST.4.0.0.0 data set with the independent ARC
1006 satellite retrievals (Section 2.5.2). The adjustments bring the *in situ* SSTs more closely in to line
1007 with the ARC SSTs. There are some discrepancies here, most notably in the early 1990s. During
1008 this period, the ARC SSTs are based on data from the ATSR1 instrument (1991-1995), which is
1009 less well understood than its successors, ATSR2 (1995-2003) and AATSR (2002-2012).

1010 Differences are larger at an individual ocean basin level and locally (Figure 12(b)).

1011 HadSST.4.0.0.0 is warmer than ARC south of 40°S and along the climatological ice edge in the
1012 Arctic. HadSST.4.0.0.0 is cooler than ARC through much of the tropics, particularly in the
1013 Indian Ocean and Maritime Continent, but differences are typically less than 0.2°C and very

1014 widely less than 0.1°C . Future versions of HadSST may benefit from using more highly-resolved
1015 adjustments. Not shown in Figure 11 is the notional 0.1°C uncertainty associated with the large-
1016 scale correlated errors in the ATSR series that arise from uncertainty in the retrieval process. It
1017 should also be noted that local differences of order $0.1\text{-}0.2^{\circ}\text{C}$ remain between ATSR and
1018 collocated buoy measurements (Merchant et al. 2012), so differences seen here are likely to be a
1019 combination of errors in ARC and errors in HadSST.4.0.0.0.

1020

1021 The comparison with Argo shows that the adjustments continue to be effective from 2012
1022 through 2017 after the failure of AATSR (Figure 11). Indeed, agreement with the independent
1023 Argo data is excellent back to 2007 when the Argo array reached its design coverage. Even
1024 before this, when coverage was less than global, the agreement remains good, suggesting that the
1025 bias adjustments are reliable at smaller scales. The map of differences between Argo and
1026 HadSST.4.0.0.0 is noisier than for the comparison with ARC due to the sparser sampling of Argo
1027 (Figure 12(c)). Some patterns are perhaps common to the two – the cooler Indian Ocean, for
1028 example – but the large differences at high latitudes are not seen in the comparison with Argo
1029 suggesting that this is due either to problems with the ARC data or is related to sampling errors
1030 of some kind.

1031

1032 Although agreement with the instrumentally homogeneous series is good overall, there are some
1033 months when the discrepancies exceed the estimated uncertainties. This is to be expected from
1034 time to time as the uncertainty range represents a 95% confidence interval. The discrepancies in
1035 the comparisons suggest an overall 1-sigma uncertainty in SST changes seen through this period
1036 of around 0.05°C , which corresponds to a stability of a few hundredths of a degree per decade. In

1037 contrast, the difference between the instrumentally homogeneous series and the unadjusted data
1038 approaches a maximum of 0.2°C (Figure 11(c)) in the global average (and locally more Figure
1039 12(d)) over the same period highlighting the importance and effectiveness of the adjustments in
1040 the modern period. The average difference between HadSST.4.0.0.0 and the instrumentally
1041 homogeneous data sets is much smaller than the applied adjustments.

1042

1043 **5.3 Comparison with other long time series**

1044 Over longer periods, it is necessary to use other data sets for comparison. We use two data sets
1045 here. The first is HadNMAT.2.0.1.0 (Kent et al. 2013) which is a data set of Nighttime Marine
1046 Air Temperatures (NMAT). Anomalies in NMAT are thought to closely track anomalies in SST
1047 over long periods and large scales (see Huang et al. 2015 for an example using climate models).
1048 The second is based on oceanographic profiles from HadIOD.1.2.0.0, excluding Argo (Atkinson
1049 et al. 2014) and adjusted using the Levitus et al. (2009) adjustments for MBTs and XBTs.

1050

1051 In order to make a direct comparison between HadSST.4.0.0.0 and HadIOD.1.2.0.0, anomalies
1052 from HadSST.4.0.0.0 were adjusted using the absolute bias rather than the relative bias so that
1053 the SSTs could be directly compared. Anomalies were then calculated for both data sets using an
1054 unadjusted climatology. However, HadSST.4.0.0.0 is not then directly comparable to
1055 HadNMAT.2.0.1.0 as HadNMAT.2.0.1.0 is provided as actuals or relative to its own adjusted
1056 1961-1990 climatology and not relative to a biased SST climatology. This will lead to a constant
1057 annual offset between the SST and NMAT series, which is approximately the size of the average
1058 climatological bias in the SST. Consequently, we shifted HadNMAT.2.0.1.0 by 0.15°C in Figure

1059 13. The offset was chosen by eye to approximately align the two series; none of the conclusions
1060 depend on the choice of offset.

1061
1062 Except for a period in the late 1940s and early 1950s, differences between HadNMAT.2.0.1.0
1063 and HadSST.4.0.0.0 anomalies (Figure 13b) on a decadal time scale are constant between 1920
1064 and 1990. Outside this period, the differences exceed the estimated uncertainties in the
1065 HadSST.4.0.0.0 data set. One notable difference occurs around 1991-1993, when
1066 HadNMAT.2.0.1.0 apparently cools relative to HadSST.4.0.0.0 (or HadSST.4.0.0.0 warms).
1067 Further investigation shows that the cooling occurs in the tropics, partly offset by warming in the
1068 northern extratropical Pacific.

1069
1070 Christy et al. (2001) previously remarked on the cooling of air temperature relative to SST in the
1071 tropics. The cause of these differential rates of warming is unknown. They hypothesized that a
1072 large scale change in circulation might have caused a persistent change in air-sea temperature
1073 differences but could not rule out the effects of biases in either the SST or NMAT data sets
1074 which, at the time, had not been studied in detail. Both HadNMAT.2.0.1.0 and HadSST.4.0.0.0
1075 are now bias adjusted and are independent of one another in the relevant period. This suggests
1076 that the differences represent a real change in the air-sea temperature difference across this
1077 transition.

1078
1079 However, other hypotheses could explain the change. There could still be an undetected bias in
1080 either the SST data, the NMAT data or both. For example, automation during the early 1990s
1081 might have allowed air temperature sensors to be placed in better-exposed locations with a

1082 consequent drop in the measured air temperature. This change would likely have been
1083 accompanied by a move to electric sensors. It could also be that the near-surface oceanographic
1084 profile data that we use as a basis for estimating the biases in ERI and modern bucket
1085 measurements change in the early 1990s in a way that is not captured by the corrections to the
1086 data. Huang et al. (2018) found a change in the average depth of profile measurements in the
1087 near-surface layer, but their criteria for selecting the profiles were different from those used here
1088 and the step change in the NMAT-SST difference is there even when unadjusted SSTs are used.

1089

1090 **5.4 Comparison with HadSST.3.1.1.0**

1091 Figure 14 shows a comparison between HadSST.4.0.0.0 and HadSST.3.1.1.0. The overall
1092 evolution of these is similar, although differences at some times are larger than the estimated
1093 uncertainties. HadSST4 runs colder than HadSST3 in the period following the Second World
1094 War to 1970. From the late 1970s to the early 2000s, HadSST4 is warmer. This change from
1095 cooler to warmer, leads to a slight sharpening of an apparent step change in global average SST
1096 around 1975. This is most distinct in the Northern Hemisphere with the transition being
1097 somewhat smoother in the Southern Hemisphere. The differences between the data sets in this
1098 period are due to two factors. First, the ERI biases are now estimated from the data (Section
1099 4.1.2) and they are larger in the 1960s in HadSST4 than they were assumed to be in HadSST3.
1100 The ERI biases are also outside the uncertainty range for ERI biases ($0.2\pm 0.1^{\circ}\text{C}$) used in
1101 HadSST3. Second, the fraction of cooler bucket measurements is lower overall in HadSST4 than
1102 in HadSST3. This is largely due to the new method of inferring the fraction of incorrectly
1103 assigned metadata (Section 4.1.4).

1104

1105 The estimated uncertainties in the global and hemispheric averages are for the most part larger in
1106 HadSST4 than HadSST3 prior to around 1970. This is due to the wider range allowed for ERI
1107 biases (Section 4.1.2) and to changes made to the bucket corrections (Section 4.1.1).

1108 **5.5 Comparison with other SST data sets**

1109 The latest versions of ERSST, COBE-SST and HadSST all now apply adjustments to the whole
1110 SST record. Figure 15 and Figure 16 show global and regional averages from HadSST.4.0.0.0,
1111 COBE-SST-2 and ERSSTv5 (with the ensemble from ERSSTv4) calculated where the data sets
1112 have common coverage. The overall evolution of the three data sets and the interannual
1113 variability in each are very similar.

1114

1115 The adjustments applied in each of the three data sets decrease the overall temperature change
1116 seen from the nineteenth century (and especially since 1900) relative to the unadjusted data. Of
1117 the three data sets, HadSST4 has a marginally higher trend from 1900 (estimated using ordinary
1118 least squares) but the difference between the trends in the three data sets is not larger than the
1119 estimated uncertainty (estimated using the ensemble with each ensemble member additionally
1120 perturbed by a sample from the measurement and sampling errors). ERSST and COBE-SST
1121 warm at a similar rate to the unadjusted data from the 1940s, 50s and 60s, but HadSST4 warms
1122 somewhat faster than the other data sets due to the adjustments applied to account for the general
1123 decline in ERI biases over that period (Figure 6). From start dates in 1970, 1980 and 1990,
1124 COBE-SST-2 warms faster than the unadjusted data and, from 1980 and 1990, faster than either
1125 HadSST4 or ERSSTv4 by a significant margin.

1126

1127 From 2000-2012, the rates of warming in all three data sets are very similar and consistent within
1128 their uncertainty ranges. All three warm faster than the unadjusted data, which has a trend close
1129 to zero. During this period, there are two important factors. First, there is a large increase in the
1130 relatively cooler drifting buoy measurements and, second, there is a decrease in the average ship
1131 bias. The analysis of HadSST4 supports ERSSTv4 and ERSSTv5 in this period (Karl et al. 2015)
1132 and is consistent with instrumentally homogeneous reference series, supporting the analysis of
1133 Hausfather et al. (2017).

1134
1135 Kent et al. (2017) showed that there were significant differences between HadSST.3.1.1.0 and
1136 ERSSTv4 at other times. The period of the largest global differences was found to be during
1137 1945-1970 when HadSST3 was warmer than ERSSTv4. HadSST4 is much closer to ERSSTv4
1138 during this period. This change is due to the new ERI bias estimates being larger than assumed in
1139 HadSST3 during this period. However, from 1960 or 1970, HadSST4 warms faster than
1140 ERSSTv4. The long-term bias adjustments in ERSSTv4 are derived from assuming a constant
1141 relationship with HadNMAT2 that, as we have already shown in Section 5.3, warms less than
1142 HadSST4 over this period with much of the difference arising from a step-like change of
1143 unknown origin in the early 1990s.

1144

1145 **6 Summary**

1146 In this paper, we have estimated biases associated with different methods for making SST
1147 measurements by comparison to near-surface oceanographic measurements and buoys. The
1148 estimated biases were combined with other metadata to bias adjust a composite SST data set.
1149 Because many of the parameters in the bias adjustment scheme are uncertain and give rise to

1150 complicated covariance structures, we present the data set as an ensemble in which we vary
1151 uncertain parameters to understand their impact on the indicators that can be derived from the
1152 data, such as the global average, or changes in temperatures.

1153

1154 The method builds on that used to create HadSST.3.1.1.0. We now have improved estimates of
1155 the biases associated with different measurement methods – including engine room
1156 measurements and insulated buckets – and we are better able to constrain poorly known
1157 parameters such as the timing of the transition from canvas to rubber buckets and the fraction of
1158 incorrect metadata on measurement method. Some of the newly constrained parameters,
1159 particularly biases associated with engine room measurements, were outside the previously
1160 estimated ranges at some times. This highlights the difficulties, expounded at greater length in
1161 Kent et al. (2017), of working with historical meteorological data, particularly when trying to
1162 make data sets that are useful for climate research. Nonetheless, by paying careful attention to
1163 the data, quantifying the biases and estimating the uncertainties, we can produce a climate data
1164 record of SST back to 1850 that is consistent with independent information.

1165

1166 The method relies on comparisons with sub-surface data. This raises two possible difficulties.
1167 First, the depth of near-surface subsurface measurements is usually slightly greater than the
1168 depth at which drifting buoys make measurements and may have changed systematically over
1169 time (Huang et al. 2018). This could lead to a cool bias in the earlier data where sub-surface
1170 measurements are used to estimate biases. On the other hand, there is no clear signal that this is
1171 the case where we can compare ship measurements to both buoys and sub-surface measurements.
1172 Second, sub-surface temperature measurements also exhibit biases. While these biases are

1173 expected to be smaller than those in SST measurements from ships, they are nonetheless
1174 significant on longer time- and space- scales. Adjustments for these biases are themselves
1175 uncertain and an active area of research (Abraham et al. 2013, Cheng et al. 2016).

1176

1177 The small adjustments that Huang et al (2015, 2017) applied to ERSSTv4/v5 in the post war
1178 years were somewhat puzzling because they suggested a small net bias during a period that saw a
1179 transition from canvas to insulated buckets and from widespread bucket use to widespread ERI
1180 use, factors that led to the larger adjustments applied to HadSST.3.1.1.0. The work we present
1181 here is in closer agreement with ERSST and suggests that the smaller net corrections are due to a
1182 greater prevalence of ERI measurements – supporting the conclusions of Carella et al. (2018) –
1183 partly offsetting larger biases associated with uninsulated buckets and an earlier change to
1184 insulated bucket use.

1185

1186 Important uncertainties likely remain. In the unadjusted data, a rapid drop in global average SST
1187 marks the end of the Second World War (Thompson et al. 2008). The drop is seen in both
1188 ERSSTv5 and HadSST.3.1.1.0, though it is less marked in the latter. It coincides with a large
1189 change in the areas sampled by the global fleet, which likely explains some of the fall. The
1190 question of how much of the remainder is artificial is still open. The separated bucket and engine
1191 room data sets considered in Section 5.1 suggest that some of the drop arises from a change in
1192 the mix of the two rather than a globally coordinated drop in actual SST. However, it is
1193 important to note that the sampling of both these data sets changes at this point. The drop is most
1194 pronounced in ERSSTv5 and coincides with a similar rapid drop seen in HadNMAT.2.0.1 at that
1195 point. However, comparisons between HadNMAT.2.0.1.0 and CRUTEM4 at common coastal

1196 grid cells (Cowtan et al. 2018) suggests that the NMAT, and hence ERSST, is artificially warm
1197 during the war years despite the adjustments that have been applied for non-standard exposure
1198 (Kent et al. 2013). These lines of evidence suggest that at least some of the drop is artificial, but
1199 they do not help to understand which of the data sets provides a better estimate. Consequently,
1200 considerable uncertainty remains regarding SST during the Second World War. This uncertainty
1201 is partly reflected in the wide uncertainty ranges given in HadSST.4.0.0.0, but a more
1202 satisfactory solution is needed. Users of the data set should be wary of drawing strong
1203 conclusions based on trends that start or end during the war years until this is resolved.

1204

1205 From 2000 to 2012, the period studied in detail here, ERSSTv5 and HadSST.4.0.0.0 have trends
1206 that are consistent with each other and with COBE-SST-2. In addition, both HadSST.4.0.0.0 and
1207 ERSSTv5 compare well with independent and instrumentally homogeneous data over the period
1208 1991 to 2017 (see also Hausfather et al. 2017). We highlight the importance of changing ship
1209 biases as well as the shifting balance of ship and buoy measurements for understanding this
1210 period.

1211

1212 Although HadSST.4.0.0.0 and ERSSTv5 show reasonable agreement in the overall evolution of
1213 global average SST, there are some interesting differences between the trends estimate from
1214 these data sets. In particular, warming since the 1950-1970 period is higher in HadSST.4.0.0.0.
1215 This is associated with a cooling of ERSSTv5 relative to HadSST.4.0.0.0 in the early 1990s. This
1216 discrepancy is also seen in a comparison with HadNMAT.2.0.1.0, the data set used to adjust
1217 ERSSTv5. The discrepancy between HadNMAT.2.0.1.0 and HadSST.4.0.0.0 suggests that there
1218 is either a large-scale change in atmospheric circulation in the early 1990s that modified the air-

1219 sea temperature difference throughout the tropics or that undetected biases remain in one or the
1220 other of the marine temperature (SST or NMAT) data sets considered here.

1221
1222 Huang et al. (2015) showed that NMAT-SST differences exhibited little variability at annual
1223 time scales in a climate model between 60°S and 60°N, although there was a long-term warming
1224 of MAT relative to SST of around 0.1°C. Hawkins et al. (2015) likewise found that MAT
1225 warmed faster globally than SST in a range of climate models. Neither paper shows specific
1226 step-change behaviour, but both show spikes in the mean model response of SST-MAT
1227 following large tropical volcanic eruptions such as Mount Pinatubo in 1991. While the step
1228 change is in the same direction as this model response, the subsequent recovery and long-term
1229 warming of MAT is not apparent in Figure 13. Hawkins et al. (2015) note that the size of the
1230 SST-MAT differences are comparable to the uncertainties in the SST and MAT data sets used,
1231 HadSST3 and HadNMAT2 in their case, but it is also true for HadSST.4.0.0.0.

1232
1233 Because of the strong links between SST and MAT and between systematic errors in SST and
1234 MAT, a fuller understanding of marine temperatures in general can only be achieved by studying
1235 both in greater detail along with metadata and other relevant marine variables such as humidity
1236 (Willett et al. 2008) and winds. While measurements of SST are now more numerous than ever
1237 thanks to the wealth of satellite data and autonomous platforms such as drifting buoys, there has
1238 been a marked continuing decline in the MAT observing system which relies on ship-borne
1239 instruments and is currently far below the level of adequacy as judged by a number of criteria
1240 (Berry and Kent 2017).

1241

1242 Recently attention has been drawn to spatial as well as temporal heterogeneity in SST biases and
1243 how these affect the interpretation of climate variability (Huang et al. 2013). Although the
1244 methods we describe account for the spatial heterogeneity that arises from geographically-
1245 varying ERI biases and changing numbers of buckets and ERI measurements, some of the factors
1246 in the adjustments – for example, the scaling for the patterns of bucket biases and how
1247 measurements without a method are assigned – are only specified at a global level and thus
1248 might not be as effective at a basin scale.

1249

1250 One factor that might vary locally is the type of bucket used for measurement. Buckets issued by
1251 different countries are of varied design and the design can affect the rate of heat loss (Carella et
1252 al. 2017b) as well as other properties of the measurements (Kent and Taylor 2006). This might
1253 be of particular importance during the early and middle decades of the 20th century when bucket
1254 use was widespread. The period 1900-1940 saw an overall increase in shipping, large changes in
1255 ship routes allowed by the opening of the Panama Canal, two World Wars, and large biases
1256 associated with uninsulated buckets deployed from fast ships. The results for modern bucket
1257 biases (Section 4.1.3) suggest a potentially important role for solar heating of the bucket and
1258 water sample. In contrast, Kent and Kaplan (2006) and Carella et al. (2017b) focused on
1259 situations in which solar heating was negligible. Recently, Chan and Huybers (2019) showed that
1260 there are relative biases between bucket measurements made by ships from different countries
1261 and between bucket measurements found in different decks in ICOADS. They argue that
1262 correcting for biases between decks and nations should improve estimated SSTs.

1263

1264 After a long analysis and discussion of the problems with the data, it can be easy to forget the
1265 enormous value that the voluntary observing ships (VOS) provide. This would be a mistake. For
1266 much of the historical record, reports from ships are all we have and, although it is the outliers of
1267 the distribution which often draw our attention, most ships have provided useful, reliable
1268 measurements. At those times where we have used other sources to correct the ship data
1269 (oceanographic profile measurements from the 1950s to present and drifting buoy measurements
1270 from the early 1990s) the ship data provide vital spatial detail in the large areas not covered by
1271 these measurements (Figure 18 and Section A2). Even in the modern period, when the coverage
1272 of drifting and moored buoy data is quasi-global, the density of shipping, particularly in the
1273 northern hemisphere, adds additional useful information (see e.g. Figure 2 and Figure 3). As
1274 researchers extend climate data records of SST derived from satellites further back in time, they
1275 will need to rely on ships to provide a “ground truth” against which their products can be tested.
1276 Last of all, ships also measure variables other than SST – air temperature (Kent et al. 2013, Berry
1277 et al. 2004), humidity (Willett et al. 2008), cloud, pressure and wind (Berry and Kent 2011) –
1278 that are essential for understanding the continual fluxes of heat and water between the
1279 atmosphere and oceans.

1280

1281 Challenges remain for building models – be they statistical or physical – which can adequately
1282 describe and constrain the spectrum or hierarchy, of errors that exist in *in situ* marine
1283 measurements of sea-surface temperature and air temperature. The methods detailed here –
1284 which can in some cases extract useful information about the error characteristics of individual
1285 ships – could be extended to include more detailed error models which track weather-dependent
1286 biases associated with ships from a particular country or which use a certain kind of bucket.

1287

1288 Finally, we reiterate the recommendations made in Kent et al. (2017), “A Call for New
1289 Approaches to Quantifying Biases in Observations of Sea Surface Temperature”, in particular
1290 the need to:

- 1291 • add more data and metadata to ICOADS;
- 1292 • reprocess existing ICOADS records;
- 1293 • improve information on observational methods;
- 1294 • improve physical and statistical models of SST bias;
- 1295 • maintain and extend the range of different estimates of SST bias; and
- 1296 • expand data sources for validation and extend the use of measures of internal consistency
1297 in validation.

1298 **Data availability**

1299 The HadSST.4.0.0.0 data set, and supporting information, is available from
1300 <http://www.metoffice.gov.uk/hadobs/hadsst4>. The following listed data sets were used in this
1301 analysis. Links to the data sets are provided where applicable.

1302

1303 International Comprehensive Ocean-Atmosphere Data Set (ICOADS) Release 3, Individual
1304 Observations. Research Data Archive at the National Center for Atmospheric Research,
1305 Computational and Information Systems Laboratory. <https://doi.org/10.5065/D6ZS2TR3>.
1306 (Research Data Archive et al. 2016)

1307 Drifting buoy data were collected and made freely available by the Copernicus project and the
1308 programs that contribute to it. Data downloaded (6 April 2018) from

1309 <http://marine.copernicus.eu/services-portfolio/access-to->

1310 [products/?option=com_csw&view=details&product_id=INSITU_GLO_NRT_OBSERVATION](#)
1311 [S_013_030](#)

1312 HadNMAT.2.0.1.0 was downloaded from <https://www.metoffice.gov.uk/hadnmat2>

1313 HadIOD is available from the corresponding author of Atkinson et al. (2014)

1314 The operational version of ERSSTv4 was downloaded from

1315 <http://www1.ncdc.noaa.gov/pub/data/cmb/ersst/v4/netcdf/>

1316 The ERSSTv4 ensemble was downloaded (16 June 2015) from

1317 <https://www1.ncdc.noaa.gov/pub/data/cmb/ersst/v4/ensemble/>

1318 The operational version of ERSSTv5 was downloaded (9 February 2018) from

1319 <https://www1.ncdc.noaa.gov/pub/data/cmb/ersst/v5/netcdf/>

1320 COBE SST 2 was provided by the NOAA/OAR/ESRL PSD, Boulder, Colorado, USA, from

1321 their Web site at <https://www.esrl.noaa.gov/psd/>

1322 (<https://www.esrl.noaa.gov/psd/data/gridded/data.cobe2.html>, accessed 23 October 2014)

1323 ARC data were downloaded (5 July 2017) from

1324 <http://catalogue.ceda.ac.uk/uuid/e6497acddf9cd8345ffbd0643c0d9729>

1325

1326 **Acknowledgments**

1327 The authors were supported by the Joint UK BEIS/Defra Met Office Hadley Centre Climate

1328 Programme (GA01101). Thanks to Boyin Huang of NOAA NCEI for making ERSSTv5 and the

1329 ERSSTv4 ensembles available and Shoji Hirahara of JMA for providing the COBE-SST-2 data

1330 set. This research would not have been possible without the vital ongoing work of the ICOADS

1331 team. This study has been conducted using E.U. Copernicus Marine Service Information. Many

1332 thanks to the reviewers for their helpful comments.

1333 **Appendices**1334 **A1 Interpolation**

1335 In order to estimate biases for different measurement types during the modern period, we used a
 1336 simple statistical interpolation scheme based on Gaussian Processes (Rasmussen and William,
 1337 2006). The scheme is described in this section. It requires as input the error covariances
 1338 calculated in Section 3.2. In addition, an estimate of the covariance matrix of the actual SST
 1339 fields is needed. The method for estimating this matrix is described in Section A1.1. The method
 1340 is applied in Section A2 and we present some tests of the method in Section A3.

1341

1342 The variability of gridded SST anomalies was modelled as a multivariate normal distribution
 1343 with mean zero and covariance matrix C . For a vector of observations y with error covariance R
 1344 (see Section 3.1 and 3.2), a globally complete reconstruction of the SST anomaly field, μ , can be
 1345 obtained using

$$\mu = CH^T(HCH^T + R)^{-1}y \quad \text{Equation A1}$$

1346 where H is a matrix consisting of zeroes and ones that selects points from C at the measured
 1347 locations in y . The posterior distribution for the reconstruction is a multivariate normal
 1348 distribution with mean, μ , and covariance P

$$P = C - CH^T(HCH^T + R)^{-1}HC \quad \text{Equation A2}$$

1349

1350 Equations A1 and A2 are equivalent to Equations 2.23 and 2.24 from Rasmussen and Williams
 1351 (2006) but with a non-diagonal error covariance. As well as producing a reconstruction of the
 1352 SST field, this framework can also be used to get an improved estimate of the errors in the data.

1353 For example, considering R as the sum of several error components $R_1, R_2 \dots R_n$ then an
 1354 improved estimate of the component R_1 is given by

$$\mu_{R_1} = R_1(HCH^T + R)^{-1}y \quad \text{Equation A3}$$

1355
 1356 with posterior covariance

$$P_{R_1} = R_1 - R_1(HCH^T + R)^{-1}R_1 \quad \text{Equation A4}$$

1357 This breakdown of the estimated errors into individual components is used in Section A3 to test
 1358 that the interpolation method is working as expected.

1359

1360 **A1.1 Estimating the prior covariance**

1361 In order to get a reasonable interpolation, a good estimate of C is needed. Here we used a time-
 1362 invariant estimate of C calculated from monthly data. C was built from a set of local covariance
 1363 kernels, in which the covariance at each point was modeled as a simple local covariance kernel.
 1364 The method is based on that used in Karspeck et al. (2012), assuming an exponential kernel and
 1365 that the length scales are strictly zonal and meridional. The covariance between two points is
 1366 equal to

$$C(x, x') = \sigma\sigma' \frac{|\Sigma|^{\frac{1}{4}}|\Sigma'|^{\frac{1}{4}}}{|\bar{\Sigma}|^{\frac{1}{2}}} \exp(-\bar{\tau}) \quad \text{Equation A5}$$

1367 Where unprimed and primed variables indicate the values at the two different points and τ , the
 1368 Mahalanobis distance, can be written as:

$$\bar{\tau}(x, x') = \sqrt{(x - x')^T \bar{\Sigma}^{-1} (x - x')} \quad \text{Equation A6}$$

1369 and,

$$\Sigma = \begin{bmatrix} L_x^2 & 0 \\ 0 & L_y^2 \end{bmatrix}; \bar{\Sigma} = \frac{1}{2}(\Sigma + \Sigma') \quad \text{Equation A7}$$

1370 Here x and y are the angular separations in longitude and latitude, L_x and L_y are the length scales.

1371 The length scales and variance of the process can vary from point to point. The primes indicate

1372 the values of Σ and σ at location x' and the unprimed variables indicate the values at x .

1373

1374 The variance was estimated first using ARC data for each grid cell by assuming that the

1375 anomalies in a grid cell were normally distributed with mean zero and variance equal to the sum

1376 of the variance of the process and the error variance (provided with the ARC data). The value of

1377 σ^2 that maximized the likelihood of the data given the parameters was found using a simple line

1378 search with increments of 0.005 K^2 . Values of the standard deviation above 1.2°C were set to

1379 1.2°C .

1380

1381 The length scales were then estimated for each grid box separately. x and y distances were

1382 measured as angular distances in latitude and longitude multiplied by 6400km with the implicit

1383 assumption that the Earth is a cylinder. This is more numerically stable than assuming a spherical

1384 Earth and the geographically-varying length scales allow for the geometrical fact that one degree

1385 of longitude corresponds to different lengths at different latitudes as well as changes in the

1386 physical length scale. Time series of anomaly differences were calculated between the target grid

1387 box and all neighbors within 10 000km and detrended using a 5th order polynomial. The true SST

1388 anomalies at the two locations were assumed to be normally distributed (mean zero and standard

1389 deviations σ and σ') and correlated with each other with a correlation of $\exp(-\tau)$. The

1390 uncertainty on the anomalies were σ_{error} and σ'_{error} . The differences between the two series are

1391 then normally distributed with mean zero and variance σ_{diff} given by

1392

$$\sigma_{diff}^2 = \sigma_{error}^2 + \sigma'_{error}{}^2 + \sigma^2 + \sigma'^2 - 2\sigma\sigma'\exp(-\tau) \quad \text{Equation A8}$$

1393

1394 The values of L_x and L_y , which maximized the likelihood of the data given the parameters were
 1395 found using the using the downhill simplex method of Nelder and Mead (1965) (as implemented
 1396 in the IDL 8.2 AMOEBA function). All other values – the error variances (σ_{error}) and the
 1397 variance of the process (σ and σ') – were fixed from earlier calculations. Missing values of L_x
 1398 were set to 2500km. Missing values of L_y were set to 1000 km. The resulting fields of σ , L_x and
 1399 L_y (Figure 17) were then stitched together to produce a single covariance matrix using equation
 1400 A5.

1401

1402 The covariances vary in character from place to place (see Figure 17). In the Tropical Pacific,
 1403 zonal length scales are long and meridional scales are limited. Variances are also higher in the
 1404 tropical Pacific associated with ENSO variability. In the North Pacific, length scales are shorter
 1405 in general. In the North Atlantic and Indian Ocean, the covariances are more isotropic with
 1406 similar zonal and meridional length scales. Over the western boundary currents, the length scales
 1407 are short and variability is high.

1408 **A2 Estimating biases**

1409 Engine Room biases and other measurement method biases were estimated using the simple
 1410 interpolation scheme in two or three steps. In the first step, gridded (see Section 3) drifting buoy
 1411 observations (if they were available) were interpolated using the formulas (Equations A9 and
 1412 A10) below to get an improved estimate of the global SST field. The mean and covariance of the
 1413 posterior distribution were used as the prior for the second step, in which near-surface sub-

1414 surface measurements (when these were available) were interpolated using the output from the
 1415 drifting-buoy interpolation as input (Equations A11 and A12). In the third step, gridded Engine
 1416 Room measurements (a similar thing can be done for other measurement types) were
 1417 interpolated using the output from the previous interpolation (Equations A13 and A14). In
 1418 contrast to the first and second, the aim in the third step was to estimate the correlated errors in
 1419 the gridded ERI measurements rather than to get an improved estimate of the SST field (although
 1420 this can be obtained as well).

1421

1422 First the buoy data are assimilated.

$$\mu_{buoy} = C_{ARC} H^T (H C_{ARC} H^T + R_{buoy})^{-1} y_{buoy} \quad \text{Equation A9}$$

$$C_{buoy} = C_{ARC} - C_{ARC} H^T (H C_{ARC} H^T + R_{buoy})^{-1} H C_{ARC} \quad \text{Equation A10}$$

1423 Where C_{ARC} is the prior covariance calculated in A1.1 using the ARC data. Then the sub-surface
 1424 data

$$\mu_{sub} = C_{buoy} H^T (H C_{buoy} H^T + R_{sub})^{-1} (y_{sub} - \mu_{buoy}) + \mu_{buoy} \quad \text{Equation A11}$$

$$C_{sub} = C_{buoy} - C_{buoy} H^T (H C_{buoy} H^T + R_{sub})^{-1} H C_{buoy} \quad \text{Equation A12}$$

1425 Finally, the bias in the ERI measurements and its uncertainty are estimated.

$$\mu_{ERIBias} = R_{ERI} (H C_{sub} H^T + R_{ERI})^{-1} (y_{ERI} - \mu_{sub}) + \mu_{sub} \quad \text{Equation A13}$$

$$R_{ERIBias} = R_{ERI} - R_{ERI} (H C_{sub} H^T + R_{ERI})^{-1} R_{ERI} \quad \text{Equation A14}$$

1426

1427 The estimated ERI bias and its uncertainty was calculated for each month using a value of σ_G of
 1428 0.2°C (the mean ERI bias used in Kilb). The resulting time series along with similar series for
 1429 buckets are shown in Figure 6. Example input and outputs for the interpolation are shown in
 1430 Figure 18a-d. Common patterns can be seen between the anomaly fields estimated using ship

1431 data (a) and buoy data (b), but there are also some differences. For example, there are features in
1432 the ship data that follow common shipping routes which could correspond to measurement
1433 errors, particularly micro biases. In the analyzed anomaly field (c), the patterns that are common
1434 to both ship and buoy data have been picked out and some of the errors that can be identified by
1435 eye in the ship data have been successfully separated (d). The largest estimated errors are in the
1436 tropics and Southern Hemisphere where ship traffic is less frequent and individual ships can have
1437 a larger effect. However, in the North Pacific and Atlantic, the ship data provide vital detail
1438 where there is no buoy coverage and the aggregate biases are less pronounced.

1439

1440 In the bias-adjustment algorithm (Section 4.1.2), a time-varying, temporally-smoothed field of
1441 μ_{ERI} was used. This ensures that the ERI measurements are unbiased relative to the drifting
1442 buoys and oceanographic data, on longer space- and time-scales, but preserves individual ship
1443 biases, which are described by the error covariances.

1444 **A3 Tests of the interpolation method**

1445 We tested the interpolation method in three different ways, which probe different aspects of the
1446 reconstruction.

1447

1448 First, we looked at its ability to reconstruct data that had been deliberately withheld (Section
1449 A3.1). Data were withheld in three different ways: at random locations, by reducing coverage in
1450 a well-observed period to match that of the 19th century and at the locations of Argo floats. These
1451 probe different aspects of the reconstruction such as the ability to fill small gaps, to reconstruct
1452 large missing areas, and to estimate SST anomalies from an independent validation system.

1453

1454 Second, we looked at the method’s skill in estimating biases for individual ships (Section A3.2).
 1455 We compared our estimates to those calculated by comparing ships to a satellite-based analysis.
 1456 Although we do not use the biases calculated for individual ships in HadSST.4.0.0.0, estimating
 1457 them tests the method’s ability to reconstruct biases locally.

1458

1459 Third, we generated realistic synthetic data for which the uncorrelated, micro bias and macro
 1460 bias errors were known (Section A3.3). We then interpolated the data and reconstructed the
 1461 micro and macro biases. This test ensures that the method can estimate the large-scale biases in
 1462 different data sources under a different set of circumstances.

1463 **A3.1 Testing the reconstruction using withheld data**

1464 In order to test the reconstruction method, fields from 2000 to 2014 were reduced in coverage.
 1465 The coverage was reduced in three ways. First, coverage in the period 2000-2014 was reduced to
 1466 that of 1850-1864. In the second test, half of the grid boxes were removed randomly. In the third
 1467 test a quarter of grid cells for which Argo measurements were available were removed. The
 1468 reduced-coverage SST fields were then reconstructed and the reconstruction was compared to the
 1469 data that had been withheld, and in the third case also to Argo data in those grid cells that had
 1470 been withheld. Two tests were then made.

1471

1472 In the first test the chi-squared statistic was calculated using,

$$\chi^2 = (\mu - w)^T (P + R)^{-1} (\mu - w) \quad \text{Equation A15}$$

1473 where μ is a vector containing the reconstructed SSTs at the locations of the withheld data, w . P
 1474 is the posterior covariance of the reconstruction and R is the error covariance of the withheld
 1475 data. For a good reconstruction – one for which the estimated fields and uncertainties are

1476 consistent – the statistic should follow a chi-squared distribution for which the number of
1477 degrees of freedom is equal to the rank of the covariance matrix, in this case equal to the number
1478 of withheld grid boxes (Povey et al. 2015).

1479

1480 The second test compared the withheld data to samples from the posterior distribution of the
1481 reconstruction, P , which were combined with samples drawn from the estimated error covariance
1482 of the withheld data, R . The combined samples give a set of fields that should resemble the
1483 withheld data – in both true SST variability and in the spectrum of observational error. The
1484 residual differences between the samples and the reconstruction were divided by the estimated
1485 uncertainty (scaled samples) as were the differences between the withheld data and the
1486 reconstruction (scaled observations). For a good reconstruction with well-specified observational
1487 uncertainties, the normalization step would, in the long run, yield distributions that are close to
1488 Normal with unit variance. However, for individual months, the spatial autocorrelation
1489 represented in P and R leads to distributions that do not look Normal.

1490

1491 The overall goodness of fit was assessed by examining how the distribution of scaled
1492 observations diverged from a distribution calculated from the scaled samples. A histogram of the
1493 scaled observations was calculated (with 0.1 unit bins) for all withheld data between 2000 and
1494 2014. An equivalent histogram of the scaled samples was calculated for the same period and then
1495 repeated 500 times with different samples. If the histogram of scaled observations falls within
1496 the range of the 500 histograms calculated from the scaled samples then the fit is considered a
1497 good one.

1498

1499 Figure 19(a) shows the distribution of the probability of the calculated chi-squared for each
1500 month between 2000 and 2016 and for each of the three different data-reduction schemes. For a
1501 good reconstruction, the probability ought to be approximately uniformly distributed. However,
1502 in this case there is a predominance of probabilities clustered towards one and zero, which
1503 suggests that the uncertainties are often underestimated and sometimes overestimated, or that the
1504 distribution of errors is not normal.

1505
1506 Figure 19 (b-d) compares the observed and theoretical distributions using the sampling method.
1507 For the case where data are “missing at random”, the observed discrepancies are most likely to
1508 fall above the range defined by the samples where deviations are small - within 0.5 standard
1509 deviations – or where deviations are very large – in excess of 3 standard deviations. To
1510 compensate, the observed discrepancies are below the expectation in most other places. Such
1511 sharp-peaked, long-tailed distributions are characteristic of observational errors in ship and buoy
1512 data (K11c) and the shape would explain the poor chi-squared distribution. Considering the
1513 simplicity of the model, the agreement between the modeled and observed distributions is rather
1514 good. In the more challenging case where the coverage is reduced to that of the 19th century, the
1515 distribution is slightly broader than the theoretical case, suggesting a slight underestimate of the
1516 uncertainty.

1517
1518 The Argo test is also encouraging. The bias in the reconstruction relative to the higher accuracy
1519 Argo measurements that arises when using the warm-biased unadjusted SST data is removed
1520 when using the bias-adjusted SSTs (Figure 19d). The distribution falls below the expected range
1521 between ± 2 standard deviations. There is a small excess at higher deviations which is associated

1522 with measurements made in the western boundary currents, marginal sea ice areas and other
1523 areas where sampling uncertainty is typically higher (assessed from maps of mean absolute
1524 residuals and root-mean-squared residuals, not shown). This suggests that either the sampling
1525 uncertainty is underestimated in these regions, or that the reconstruction technique cannot
1526 resolve the small scale variability in these regions.

1527

1528 An additional test was done based on the Argo data. Data were removed at the locations of the
1529 selected Argo data, but instead of comparing the reconstruction to the Argo data, it was
1530 compared to the withheld SST data. The shape of the distribution is somewhere between that
1531 seen in the missing-at-random case and that seen in the 19th Century coverage case. Given the
1532 good agreement between Argo and the reconstruction, the implication is that the narrow-peak,
1533 long-tail shape of the distribution of differences between SST and reconstruction arises purely
1534 from the distribution of observational errors in the SST data.

1535

1536 K11c investigated the effect of errors that are not normally distributed and which vary from
1537 agent to agent. They derived representative uncertainties that provide a good overall fit to the
1538 spectrum of errors seen in real data and it is these that we use in our analysis. Some care is
1539 needed when interpreting the standard deviations obtained from the data in terms of an actual
1540 distribution of errors because they cannot be assumed to be Normal. However, these tests
1541 indicate that there are no severe biases in the method.

1542

1543 **A3.2 Testing the reconstruction by estimating individual ship biases**

1544 It is possible to extract a posterior distribution for the error characteristics of each identifiable
 1545 ship, which can then be compared to other estimates as a way of checking the reliability of the
 1546 reconstruction method. This is done by noting that the error covariance, R , in the above equations
 1547 is a simple sum of contributions from individual ships. Each of these individual contributions can
 1548 be used to get an improved mean and covariance for the error characteristics of individual ships
 1549 and buoys. Of particular interest are the micro biases, which are an important component of the
 1550 measurement uncertainty at a global level.

1551

1552 It is worth pausing here momentarily to remind ourselves of the discussion and definitions of
 1553 “error” and “uncertainty” from Section 2.1 and Section 3 because the nomenclature can get
 1554 horribly confusing at this point. What we are attempting to estimate are the values of B_{ij} from
 1555 Equation 3.0. B_{ij} is that part of the error (defined as the difference between the true SST and the
 1556 measured SST) that is a persistent offset associated with a particular ship and for simplicity’s
 1557 sake we shall assume that B_{ij} is constant for a particular ship, i , so we can write B_i . To start with
 1558 we assume that B_i is zero for all ships, with a large uncertainty, σ_{b_i} . Using the method described
 1559 below, we can make an improved estimate of the size of the error, B'_i , and its uncertainty

1560 $\sigma'_{b_i} < \sigma_{b_i}$.

1561

1562 For each ship, i , the individual contribution, R_i to the overall error covariance (individual terms
 1563 in the summation in Equation 3.4) was estimated and then the posterior mean and covariance of
 1564 the errors for that ship were estimated using.

$$\mu_i = R_i(HCH^T + R)^{-1}y \quad \text{Equation A16}$$

$$R'_i = R_i - R_i(HCH^T + R)^{-1}R_i \quad \text{Equation A17}$$

1565 The mean, μ_i , found in this way is the weighted contribution of a particular error to the gridded
 1566 averages, so it needs to be divided by the appropriate weight (w_i , Equation 3.2) to obtain an
 1567 updated estimate of the bias, B'_i ,

$$B'_i = \frac{\mu_i}{w_i}$$

1568 Likewise the covariance R'_i can be processed to obtain an updated estimate of σ'_{b_i} .

1569

1570 We applied Equations A16 and A17 for every month from 2000 to 2012 with: y being the
 1571 median adjusted HadSST.4.0.0.0 for that month; the large-scale bias covariances in R set to zero;
 1572 and C as derived in Section A1.1. We then extracted an estimate of B_i for every uniquely
 1573 identifiable ship.

1574

1575 For ships in well-travelled regions, the resulting micro biases, B'_i , remain close to the prior: the
 1576 mean is zero and the uncertainty almost equal to σ_{b_i} . This is expected because many ships are
 1577 averaged together in well-travelled grid boxes making the individual contributions impossible to
 1578 separate. However, if a ship visits several, poorly-populated grid cells, the estimated micro bias
 1579 will take some more-definite value and the posterior variance will be small.

1580

1581 Figure 20 shows the monthly estimated ship micro biases for those ships where the estimated
 1582 uncertainty in the micro bias was less than 0.25°C compared to estimated micro biases taken
 1583 from the IQUAM (in situ QUALity Monitoring, Xu and Ignatov 2010) tool. In the IQUAM
 1584 analysis, *in situ* data were compared to a daily background SST field derived from a combination
 1585 of satellite and in *in situ* data. The correlation between the estimates from IQUAM and the

1586 HadSST.4.0.0.0 analysis is around 0.7 where the uncertainty in the HadSST.4.0.0.0-estimated
1587 bias is less than 0.25°C. Differences are expected as IQUAM uses a different QC system and
1588 therefore the micro biases are not estimated from exactly the same observations. The figure
1589 shows the regression of the IQUAM estimates on HadSST.4.0.0.0 estimates and vice versa
1590 (diagonal blue lines). The two regression lines encompass a line that is parallel to $y=x$ (the red
1591 line) and passes through 0.19°C, which can be interpreted as the approximate average bias
1592 between the two analyses during this period: the HadSST.4.0.0.0 estimate is adjusted to remove
1593 large-scale biases, but this is not done for the IQUAM data. A rough estimate of the large-scale
1594 average ship bias in HadSST.4.0.0.0 for this period can be calculated from the combined ship
1595 and buoy biases in the three regions shown in Figure 8 by dividing the bias by the fractional
1596 contribution of ships to the average. This gives a range of values from 0.14 (in the Northern
1597 Hemisphere) to 0.20°C (in the Southern Hemisphere) which encompasses 0.19°C. The analysis
1598 suggests that the magnitudes of the systematic errors for individual ships have not been
1599 systematically underestimated and that the reconstruction can reliably estimate the size of
1600 systematic errors in the data even for the challenging case of individual ship micro biases.

1601

1602 **A3.3 Testing the reconstruction using synthetic data**

1603 In the final test, a set of synthetic observations were generated from a globally-complete high-
1604 resolution (1/20° grid resolution) daily SST data set, called OSTIA (Operational Sea Surface
1605 Temperature and Sea Ice Analysis, Donlon et al. 2012). The idea is to create a synthetic, but
1606 realistic data set with known SSTs and known measurement errors. The techniques developed in
1607 this paper can then be applied to estimate the errors and compare them to the prescribed values.

1608

1609 The OSTIA SST fields were sampled at the locations of ICOADS observations and synthetic
1610 measurement errors were added to the data. If no OSTIA SST was available, no synthetic
1611 observation was produced. The errors were drawn from the error model described in Section 3.1.
1612 Each synthetic observation had an uncorrelated measurement error added to it which was drawn
1613 from a normal distribution with a standard deviation equal to the estimated uncorrelated error
1614 uncertainty for that platform (ship, drifting buoy or moored buoy, see Table 1). In addition, each
1615 individual agent, identified by its ID in ICOADS, was assigned a micro-bias error that was added
1616 to every observation made by the agent. The micro-bias errors were drawn from a normal
1617 distribution with a standard deviation equal to the estimated micro-bias error uncertainty for that
1618 platform (Table 1). Finally, characteristic biases were added to all observations made using a
1619 particular measurement method. The characteristic biases were drawn from a normal distribution
1620 with standard deviation of 0.2°C . New error values were drawn for each month for all
1621 components.

1622

1623 The synthetic observations were processed in the same way as the actual observations. The
1624 synthetic observations were gridded and the uncertainties in the gridded data were calculated
1625 (Section 3). The large-scale biases (Section A2) and micro-biases for individual ships were
1626 estimated (Section A3.2) for each month.

1627

1628 Figure 21 shows the difference between the estimated and prescribed synthetic biases for each
1629 measurement type and the estimated uncertainty envelope. The uncertainties in the biases are
1630 reasonably well estimated for the bucket and hull sensor measurements. However, there is some
1631 evidence of a slight bias in the estimation of the ERI biases, which appears to be consistently
1632 “warm” prior to 2009. This is odd, because the input biases have a mean of zero and are

1633 symmetrically distributed, and suggests that residual errors of around 0.05°C cannot be reliably
 1634 eliminated using this method. Consequently, a lower limit on uncertainties estimated in this way
 1635 is set at 0.05°C.

1636

1637 Figure 22 shows the comparison of the assigned and estimated micro-biases for ships. For the
 1638 majority of ships, the estimated micro-biases are close to zero (Figure 22(a)), which is the mean
 1639 of the prior estimate for the micro-biases. This happens because the estimates are based on
 1640 coarsely gridded data and there is insufficient information to estimate the micro-biases if many
 1641 ships contribute to the same grid-box average. This is reflected in the uncertainties attached to
 1642 each estimate of the micro-bias. Selecting only those ships where the uncertainty is significantly
 1643 lower than the prior value (Figure 22(b)-(d)) shows a closer correlation between the estimated
 1644 and assigned micro-biases. This further demonstrates the ability of the method to extract
 1645 individual ship biases.

1646 **A4 Derivation of bucket biases**

1647 The bias for a grid cell can be written (Equation 4.1) as

$$\mathbf{B} = f_e \mathbf{E} + f_c \mathbf{B}_{tc} + f_w \mathbf{B}_{tw} + f_r \mathbf{B}_{tr} + f_d \mathbf{D} \quad \text{Equation A18}$$

1648 The bucket correction for a canvas bucket, B_c , which adjusts the grid-box average bias to be
 1649 consistent with the climatological average, can be written as

$$\mathbf{B}_c = \mathbf{B}_{tc} - \bar{\mathbf{B}} \quad \text{Equation A19}$$

1650 Where the overbar denotes the 1961-1990 average. Expanding this out:

$$\mathbf{B}_c = \mathbf{B}_{tc} - \overline{f_e \mathbf{E}} - \overline{f_c \mathbf{B}_{tc}} - \overline{f_w \mathbf{B}_{tw}} - \overline{f_r \mathbf{B}_{tr}} - \overline{f_d \mathbf{D}} \quad \text{Equation A20}$$

1651 The overbars cover only those terms of the equation that are time-varying. In addition, f_w is zero
 1652 during the climatology period (there are considered to be no wooden buckets) and D , the drifter
 1653 bias, is zero at all times, so

$$B_c = B_{tc} - \overline{f_e E} - \overline{f_c B_{tc}} - \overline{f_r B_{tr}} \quad \text{Equation A21}$$

1654 And, rearranging:

$$B_{tc} = \frac{B_c + \overline{f_e E} + \overline{f_r B_{tr}}}{(1 - \overline{f_c})} \quad \text{Equation A22}$$

1655 We get the formula for calculating the true bias of a canvas bucket given the bucket correction
 1656 for a canvas bucket and the true biases for ERI and rubber bucket measurements.

1657 **A5 Comparison to AR5**

1658 The IPCC AR5 (Hartman et al. 2013) showed trends over particular periods for a number of
 1659 different SST data sets: HadISST1.1 (Rayner et al. 2003), HadSST2 (R06), COBE-SST (Ishii et
 1660 al. 2005) and ERSSTv3 (Smith et al. 2008). Figure 23 shows the HadSST.4.0.0.0 and ERSSTv4
 1661 ensemble estimates of the trends over the same periods (1880-2012, 1901-2012, 1951-2012 and
 1662 1979-2012). Also shown are trends over the period 1998-2012 (used in Karl et al. 2015) and
 1663 2002-2012.

1664

1665 There are a number of interesting things to note. First, the spread from the HadSST.4.0.0.0
 1666 ensemble, which incorporates uncertainty in the bias adjustments only, is a major part of the
 1667 uncertainty at long time scales (>30 years, Figure 23 (a)-(d)). At shorter time scales, however,
 1668 the ensemble spread contributes only a small part to the overall uncertainty, the remainder
 1669 coming from measurement and sampling errors, particularly the effect of micro-biases. This is
 1670 particularly clear for the periods 1998-2012 and 2002-2012 where this component dominates

1671 (Figure 23 (e) and (f)). Second, the ERSSTv4 and HadSST.4.0.0.0 ensembles overlap at all
1672 times, implying that the two estimates are consistent for this measure. Third, over the largest
1673 time periods (Figure 23(a) and (b)) the net effect of the adjustments is to reduce the trend relative
1674 to the unadjusted observations.

1675

1676 In the period 1998-2012, the unadjusted data have a trend close to zero. In the period 2002-2012,
1677 the unadjusted trend is negative. In contrast, all the adjusted data sets indicate more warming (or
1678 less cooling) than in the unadjusted data. This is consistent across the three data sets. For other
1679 periods, 1979-2012, for example, the sign of the correction is not clear with some adjustments
1680 increasing the trend and others reducing it even within one ensemble.

1681

1682 Finally, the estimates from the three current, fully-adjusted, data sets all tend to sit at the upper
1683 end of the range from the data sets employed in IPCC AR5.

1684 **References**

1685 Abraham, J.P., et al. (2013), A review of global ocean temperature observations: Implications for
1686 ocean heat content estimates and climate change, *Rev. Geophys.*, 51, 450-483,
1687 doi:10.1002/rog.20022.

1688

1689 Atkinson C.P., Rayner, N.A., Kennedy, J.J. and Good S.A. (2014), An integrated database of
1690 ocean temperature and salinity observations, *Journal of Geophysical Research: Oceans*, 119
1691 (2014), no. 10, 7139-7163.

1692

- 1693 Berry, D.I., Kent, E.C. and Taylor, P.K. (2004), An Analytical Model of Heating Errors in
1694 Marine Air Temperatures from Ships. *J. Atmos. Oceanic Technol.*, 21, 1198–1215,
1695 [https://doi.org/10.1175/1520-0426\(2004\)021<1198:AAMOHE>2.0.CO;2](https://doi.org/10.1175/1520-0426(2004)021<1198:AAMOHE>2.0.CO;2)
1696
- 1697 Berry, D. I. and Kent, E. C. (2017), Assessing the health of the in situ global surface marine
1698 climate observing system. *Int. J. Climatol.*, 37: 2248-2259. doi:10.1002/joc.4914
1699
- 1700 Berry, D. I. and Kent, E. C. (2011), Air–Sea fluxes from ICOADS: the construction of a new
1701 gridded dataset with uncertainty estimates. *Int. J. Climatol.*, 31: 987-1001. doi:10.1002/joc.2059
1702
- 1703 Berry, D.I.; Corlett, G.K.; Embury, O.; Merchant, C.J. (2018), Stability Assessment of the
1704 (A)ATSR Sea Surface Temperature Climate Dataset from the European Space Agency Climate
1705 Change Initiative. *Remote Sens.*, 10, 126.
1706
- 1707 Bojinski, S., M. Verstraete, T.C. Peterson, C. Richter, A. Simmons, and M. Zemp, (2014), The
1708 Concept of Essential Climate Variables in Support of Climate Research, Applications, and
1709 Policy. *Bull. Amer. Meteor. Soc.*, 95, 1431–1443, <https://doi.org/10.1175/BAMS-D-13-00047.1>
1710
- 1711 Brönnimann, S. (2005), The global climate anomaly 1940–1942. *Weather*, 60: 336-342.
1712 doi:10.1256/wea.248.04
1713
- 1714 Carella, G., Kent, E.C. and Berry, D.I. (2017a), A probabilistic approach to ship voyage
1715 reconstruction in ICOADS. *Int. J. Climatol.*, 37: 2233-2247. Doi:10.1002/joc.4492

1716

1717 Carella, G., Morris, A. K. R., Pascal, R. W., Yelland, M. J., Berry, D. I., Morak-Bozzo, S.,
1718 Merchant, C. J. and Kent, E. C. (2017b), Measurements and models of the temperature change of
1719 water samples in sea-surface temperature buckets. *Q.J.R. Meteorol. Soc.*, 143: 2198–2209.
1720 doi:10.1002/qj.3078

1721

1722 Carella, G., Kennedy, J. J., Berry, D. I., Hirahara, S., Merchant, C. J., Morak-Bozzo, S., & Kent,
1723 E. C. (2018), Estimating sea surface temperature measurement methods using characteristic
1724 differences in the diurnal cycle. *Geophysical Research Letters*, 45, 363–371.
1725 <https://doi.org/10.1002/2017GL076475>

1726

1727 Carton, J.A. and Giese, B.S. (2008), A Reanalysis of Ocean Climate Using Simple Ocean Data
1728 Assimilation (SODA). *Mon. Wea. Rev.*, 136, 2999–3017,
1729 <https://doi.org/10.1175/2007MWR1978.1>

1730

1731 Chan, D. and Huybers, P., (2019) Systematic differences in bucket sea surface temperature
1732 measurements amongst nations identified using a linear-mixed-effect method. *J. Climate*,
1733 <https://doi.org/10.1175/JCLI-D-18-0562.1>

1734

1735 Cheng, L., Abraham, J., Goni, G., Boyer, T., Wijffels, S., Cowley, R., Gouretski, V., Reseghetti,
1736 F., Kizu, S., Dong, S., Bringas, F., Goes, M., Houpert, L., Sprintall, J., and Zhu, J. (2016), XBT
1737 Science: Assessment of Instrumental Biases and Errors. *Bull. Amer. Meteor. Soc.*, 97, 924–933,
1738 <https://doi.org/10.1175/BAMS-D-15-00031.1>

1739

1740 Christy J.R., Parker, D.E., Brown, S.J., Macadam, I., Stendel, M., and Norris, W.B. (2001),
1741 Differential trends in tropical sea surface and atmospheric temperatures since 1979, *Geophysical*
1742 *Research Letters* 28(1), 183-186.

1743

1744 Clark I. (1979), *Practical Geostatistics*, London (Applied Sci. Publ. Ltd.) ISBN-0-853340843-X

1745

1746 Compo, G. P., Whitaker, J. S., Sardeshmukh, P. D., Matsui, N. , Allan, R. J., Yin, X. , Gleason,
1747 B. E., Vose, R. S., Rutledge, G. , Bessemoulin, P. , Brönnimann, S. , Brunet, M. , Crouthamel, R.
1748 I., Grant, A. N., Groisman, P. Y., Jones, P. D., Kruk, M. C., Kruger, A. C., Marshall, G. J.,
1749 Maugeri, M. , Mok, H. Y., Nordli, Ø. , Ross, T. F., Trigo, R. M., Wang, X. L., Woodruff, S. D.
1750 and Worley, S. J. (2011), The Twentieth Century Reanalysis Project. *Q.J.R. Meteorol. Soc.*, 137:
1751 1-28. doi:10.1002/qj.776

1752

1753 Cowley R., Wijffels, S., Cheng, L., Boyer, T., and Kizu, S. (2013), Biases in expendable
1754 bathythermograph data: A new view based on historical side-by-side comparisons, *Journal of*
1755 *Atmospheric and Oceanic Technology*, 30(6), 1195-1225.

1756

1757 Cowtan, K., Hausfather, Z., Hawkins, E., Jacobs, P., Mann, M.E., Miller, S.K., Steinman, B.A.,
1758 Stolpe, M. B., and Way, R. G. (2015), Robust comparison of climate models with observations
1759 using blended land air and ocean sea surface temperatures, *Geophys. Res. Lett.*, 42, 6526–6534,
1760 doi:10.1002/2015GL064888.

1761

1762 Cowtan, K., Rohde, R., Hausfather, Z. (2018), Evaluating biases in sea surface temperature
1763 records using coastal weather stations. *Q J R Meteorol Soc.* 144: 670–681.

1764 <https://doi.org/10.1002/qj.3235>

1765

1766 Donlon, C., et al. (2007), The Global Ocean Data Assimilation Experiment High-resolution Sea
1767 Surface Temperature Pilot Project. *Bulletin of the American Meteorological Society*, 88,
1768 11971213. doi: 10.1175/BAMS-88-8-1197

1769

1770 Donlon, C. J., Martin, M., Stark, J., Roberts-Jones, J., Fiedler, E. and Wimmer, W. (2012), The
1771 operational sea surface temperature and sea ice analysis (OSTIA) system. *Remote Sensing of*
1772 *Environment*, 116, 140-158.

1773

1774 Embury O., Merchant, C.J., and Corlett, G.K. (2012), A reprocessing for climate of sea surface
1775 temperature from the along-track scanning radiometers: Initial validation, accounting for skin
1776 and diurnal variability effects, *Remote Sensing of Environment* 116, 62-78, Advanced Along
1777 Track Scanning Radiometer(AATSR) Special Issue.

1778

1779 Embury, O. (2012), ARC: Level 3 daily sea surface temperature data v1.1. NERC Earth
1780 Observation Data Centre, downloaded 5 Jul 2017.

1781 <http://catalogue.ceda.ac.uk/uuid/e6497acddf9cd8345ffbd0643c0d9729>

1782

1783 Folland, C.K. and Parker, D.E. (1995), Correction of instrumental biases in historical sea surface
1784 temperature data. *Q.J.R. Meteorol. Soc.*, 121: 319-367. doi:10.1002/qj.49712152206

1785

1786 Freeman, E. , Woodruff, S. D., Worley, S. J., Lubker, S. J., Kent, E. C., Angel, W. E., Berry, D.
1787 I., Brohan, P. , Eastman, R. , Gates, L. , Gloeden, W. , Ji, Z. , Lawrimore, J. , Rayner, N. A.,
1788 Rosenhagen, G. and Smith, S. R. (2017), ICOADS Release 3.0: a major update to the historical
1789 marine climate record. *Int. J. Climatol.*, 37: 2211-2232. doi:10.1002/joc.4775

1790

1791 Garrett, J.F. (1980), Availability of the FGGE drifting buoy system data set, *Deep Sea Research*
1792 *Part A. Oceanographic Research Papers* 27(12), 1083-1086 <https://doi.org/10.1016/0198->
1793 [0149\(80\)90067-9](https://doi.org/10.1016/0198-0149(80)90067-9)

1794

1795 Good, S.A., Martin, M.J. and Rayner, N.A. (2013), EN4: quality controlled ocean temperature
1796 and salinity profiles and monthly objective analyses with uncertainty estimates, *Journal of*
1797 *Geophysical Research: Oceans*, 118, 6704-6716, doi:10.1002/2013JC009067

1798

1799 Gouretski, V. (2012), Using GEBCO digital bathymetry to infer depth biases in the XBT data,
1800 *Deep Sea Research Part I: Oceanographic Research Papers* 62(C), 40-52.

1801

1802 Gouretski V., Kennedy, J., Boyer, T., and Köhl, A. (2012), Consistent near-surface ocean
1803 warming since 1900 in two largely independent observing networks, *Geophysical Research*
1804 *Letters* 39, L19606, doi: 10.1029/2012GL052975.

1805

1806 Gouretski, V. and Reseghetti, F. (2010), On depth and temperature biases in bathythermograph
1807 data: Development of a new correction scheme based on analysis of a global ocean database,

1808 *Deep Sea Research Part I: Oceanographic Research Papers* 57(6), 812 - 833.

1809 <https://doi.org/10.1016/j.dsr.2010.03.011>

1810

1811 Hartmann, D.L., Klein Tank, A.M.G., Rusticucci, M., Alexander, L.V., Brönnimann, S.,

1812 Charabi, Y., Dentener, F.J., Dlugokencky, E.J., Easterling, D.R., Kaplan, A., Soden, B.J.,

1813 Thorne, P.W., Wild, M. and Zhai, P.M. (2013), Observations: Atmosphere and Surface. In:

1814 Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth

1815 Assessment Report of the Intergovernmental Panel on Climate Change [Stocker, T.F., D. Qin,

1816 G.-K. Plattner, M. Tignor, S.K. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex and P.M. Midgley

1817 (eds.)]. Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA.

1818

1819 Hausfather Z., Cowtan, K., Clarke, D.C., Jacobs, P., Richardson, M., and Rohde, R. (2017),

1820 Assessing recent warming using instrumentally homogeneous sea surface temperature records ,

1821 *Science Advances* 3(1), e1601207, doi: 10.1126/sciadv.1601207

1822 Hersbach, H., Peubey, C., Simmons, A., Berrisford, P., Poli, P. and Dee, D. (2015), ERA-20CM:

1823 a twentieth-century atmospheric model ensemble. *Q.J.R. Meteorol. Soc.*, 141: 2350–2375.

1824 doi:10.1002/qj.2528

1825

1826 Hirahara S., Ishii, M., and Fukuda, Y. (2014), Centennial-scale sea surface temperature analysis

1827 and its uncertainty, *Journal of Climate* 27(1), 57-75. <https://doi.org/10.1175/JCLI-D-12-00837.1>

1828

1829 Huang, B., L’Heureux, M., Lawrimore, J., Liu, C., Zhang, H., Banzon, V., Hu, Z. and Kumar, A.

1830 (2013), Why Did Large Differences Arise in the Sea Surface Temperature Datasets across the

1831 Tropical Pacific during 2012?. *J. Atmos. Oceanic Technol.*, 30, 2944–2953,

1832 <https://doi.org/10.1175/JTECH-D-13-00034.1>

1833

1834 Huang B., Banzon, V.F., Freeman, E., Lawrimore, J., Liu, W., Peterson, T.C., Smith, T.M.,

1835 Thorne, P.W., Woodruff, S.D. and Zhang, H-Min (2015), Extended reconstructed sea surface

1836 temperature version 4 (ERSST v4). Part I: Upgrades and intercomparisons, *Journal of Climate*

1837 28(3), 911-930, <https://doi.org/10.1175/JCLI-D-14-00006.1>

1838

1839 Huang B., Thorne, P.W., Smith, T.M., Liu, W., Lawrimore, J., Banzon, V.F., Zhang, H-Min,

1840 Peterson, T.C. and Menne, M. (2016), Further exploring and quantifying uncertainties for

1841 extended reconstructed sea surface temperature (ERSST) version 4 (v4), *Journal of Climate*

1842 29(9), 3119-3142, <https://doi.org/10.1175/JCLI-D-14-00007.1>

1843

1844 Huang B., Thorne, P.W., Banzon, V.F., Boyer, T., Chepurin, G., Lawrimore, J.H., Menne, M.J.,

1845 Smith, T.M., Vose, R.S., and Zhang, H-Min (2017), Extended reconstructed sea surface

1846 temperature version 5 (ERSSTv5): Upgrades, validations, and intercomparisons, *Journal of*

1847 *Climate*, 30, 8179–8205, <https://doi.org/10.1175/JCLI-D-16-0836.1>

1848

1849 Huang, B., Angel, W., Boyer, T., Cheng, L., Chepurin, G., Freeman, E., Liu, C., and Zhang, H.

1850 (2018), Evaluating SST Analyses with Independent Ocean Profile Observations. *J. Climate*, 31,

1851 5015–5030, <https://doi.org/10.1175/JCLI-D-17-0824.1>

1852

1853 Ishii, M. , Shouji, A. , Sugimoto, S. and Matsumoto, T. (2005), Objective analyses of sea-surface
1854 temperature and marine meteorological variables for the 20th century using ICOADS and the
1855 Kobe Collection. *Int. J. Climatol.*, 25: 865-879. doi:10.1002/joc.1169

1856

1857 JCGM 100 – Evaluation of measurement data – Guide to the expression of uncertainty in
1858 measurement (ISO/IEC Guide 98-3)

1859 <https://www.iso.org/sites/JCGM/GUM/JCGM100/C045315e->
1860 [html/C045315e.html?csnumber=50461](https://www.iso.org/sites/JCGM/GUM/JCGM100/C045315e.html?csnumber=50461)

1861

1862 Kalnay, E., Kanamitsu, M., Kistler, R., Collins, W., Deaven, D., Gandin, L., Iredell, M., Saha,
1863 S., White, G., Woollen, J., Zhu, Y., Leetmaa, A., Reynolds, R., Chelliah, M., Ebisuzaki, W.,
1864 Higgins, W., Janowiak, J., Mo, K.C., Ropelewski, C., Wang, J., Jenne, R., and Joseph, D. (1996),
1865 The NCEP/NCAR 40-Year Reanalysis Project. *Bull. Amer. Meteor. Soc.*, 77, 437–471,

1866 [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)

1867

1868 Karl, T.R., Arguez, A., Huang, B., Lawrimore, J.H., McMahon, J.R., Menne, M.J., Peterson,
1869 T.C., Vose, R.S. and Zhang, H-Min (2015), Possible artifacts of data biases in the recent global
1870 surface warming hiatus, *Science*, aaa5632. DOI: 10.1126/science.aaa5632

1871

1872 Karspeck, A.R., Kaplan, A. and Sain, S.R. (2012), Bayesian modelling and ensemble
1873 reconstruction of mid-scale spatial variability in North Atlantic sea-surface temperatures for
1874 1850-2008, *Quarterly Journal of the Royal Meteorological Society* 138(662), 234-248.

1875 doi:10.1002/qj.900

1876

1877 Kawai, Y. and Wada, A. (2007), Diurnal sea surface temperature variation and its impact on the
1878 atmosphere and ocean: a review. *Journal of Oceanography*, 2007, 63(5), 721-744.

1879 <https://doi.org/10.1007/s10872-007-0063-0>

1880

1881 Kennedy J.J., Rayner, N.A., Smith, R.O., Saunby, M. and Parker, D.E. (2011a), Reassessing
1882 biases and other uncertainties in sea-surface temperature observations measured in situ since
1883 1850, Part 1: measurement and sampling uncertainties, *J. Geophys. Res. Atmospheres.*, 116,
1884 D14103, doi: 10.1029/2010JD015218.

1885

1886 Kennedy J.J., Rayner, N.A., Smith, R.O., Saunby, M. and Parker, D.E. (2011b), Reassessing
1887 biases and other uncertainties in sea-surface temperature observations measured in situ
1888 since 1850, Part 2: biases and homogenisation, *J. Geophys. Res. Atmos.*, 116, D14104, doi:
1889 10.1029/2010JD015220..

1890

1891 Kennedy J.J., Smith, R.O. and Rayner, N.A. (2011c), Using AATSR data to assess the quality of
1892 in situ sea surface temperature observations for climate studies, *Remote Sensing of Environment*,
1893 116(15), 79-92, <https://doi.org/10.1016/j.rse.2010.11.021>

1894

1895 Kennedy, J.J. (2014), A review of uncertainty in in situ measurements and data sets of sea
1896 surface temperature, *Rev. Geophys.*, 52, 1-32, doi:10.1002/2013RG000434

1897

- 1898 Kent, E.C. and Taylor, P.K. (2006), Toward Estimating Climatic Trends in SST. Part I: Methods
1899 of Measurement. *J. Atmos. Oceanic Technol.*, 23, 464–475,
1900 <https://doi.org/10.1175/JTECH1843.1>
1901
- 1902 Kent, E.C. and Kaplan, A. (2006), Toward Estimating Climatic Trends in SST. Part III:
1903 Systematic Biases. *Journal of Atmospheric and Oceanic Technology*, 23:3, 487-500,
1904 <https://doi.org/10.1175/JTECH1845.1>
1905
- 1906 Kent, E.C., Kennedy, J.J., Berry, D.I. and Smith, R.O. (2010), Effects of instrumentation changes
1907 on sea surface temperature measured in situ, *Wiley Interdisciplinary Reviews: Climate Change*,
1908 1(5), 718-728, <https://doi.org/10.1002/wcc.55>
1909
- 1910 Kent, E.C., Woodruff, S.D. and Berry, D.I. (2007), Metadata from WMO publication no. 47 and
1911 an assessment of voluntary observing ship observation heights in ICOADS, *Journal of*
1912 *Atmospheric and Oceanic Technology*, 24(2), 214-234, <https://doi.org/10.1175/JTECH1949.1>
1913
- 1914 Kent, E.C., Kennedy, J.J., Smith, T.M., Hirahara, S., Huang, B., Kaplan, A., Parker, D.E.,
1915 Atkinson, C.P., Berry, D.I., Carella, G., Fukuda, Y., Ishii, M., Jones, P.D., Lindgren, F.,
1916 Merchant, C.J., Morak-Bozzo, S., Rayner, N.A., Venema, V., Yasui, S., and Zhang, H. (2017), A
1917 Call for New Approaches to Quantifying Biases in Observations of Sea Surface Temperature.
1918 *Bull. Amer. Meteor. Soc.*, 98, 1601–1616, <https://doi.org/10.1175/BAMS-D-15-00251.1>
1919

- 1920 Kent, E.C., Rayner, N.A., Berry, D.I., Saunby, M., Moat, B.I., Kennedy, J.J., and Parker, D.E.
1921 (2013), Global analysis of night marine air temperature and its uncertainty since 1880: The
1922 HadNMAT2 data set, *Journal of Geophysical Research: Atmospheres*, 118(3), 1281-1298,
1923 <https://doi.org/10.1002/jgrd.50152>
1924
- 1925 Lean, K. and Saunders, R.W. (2013), Validation of the ATSR reprocessing for climate (ARC)
1926 dataset using data from drifting buoys and a three-way error analysis, *Journal of Climate* 26(13),
1927 4758-4772, <https://doi.org/10.1175/JCLI-D-12-00206.1>
1928
- 1929 Levitus S., Antonov, J.I., Boyer, T.P., Locarnini, R.A., Garcia, H.E., and Mishonov, A.V.
1930 (2009), Global ocean heat content 1955-2007 in light of recently revealed instrumentation
1931 problems, *Geophys. Res. Lett.*, 36, L07608. <https://doi.org/10.1029/2008GL037155>
1932
- 1933 Liu W., Huang, B., Thorne, P.W., Banzon, V.F., Zhang, H-Min, Freeman, E., Lawrimore, J.,
1934 Peterson, T.C., Smith, T.M. and Woodruff, S.D. (2015), Extended reconstructed sea surface
1935 temperature version 4 (ERSSTv4): Part II. Parametric and structural uncertainty estimations,
1936 *Journal of Climate*, 28(3), 931-951, <https://doi.org/10.1175/JCLI-D-14-00007.1>
1937
- 1938 Medhaug I., Stolpe, M.B., Fischer, E.M. and Knutti, R. (2017), Reconciling controversies about
1939 the global warming hiatus, *Nature* 545(7652), 41-47.
1940
- 1941 Merchant C.J., Embury, O., Rayner, N.A., Berry, D.I., Corlett, G.K., Lean, K., Veal, K.L., Kent,
1942 E.C., Llewellyn-Jones, D.T. and Remedios, J.J. (2012), A 20 year independent record of sea

1943 surface temperature for climate from along-track scanning radiometers, *Journal of Geophysical*
1944 *Research: Oceans*, (1978-2012) 117(C12), <https://doi.org/10.1029/2012JC008400>

1945

1946 Merchant, C. J., Embury, O., Roberts-Jones, J., Fiedler, E., Bulgin, C. E., Corlett, G. K., Good,
1947 S., McLaren, A., Rayner, N., Morak-Bozzo, S. and Donlon, C. (2014), Sea surface temperature
1948 datasets for climate applications from Phase 1 of the European Space Agency Climate Change
1949 Initiative (SST CCI). *Geosci. Data J.*, 1, 179–191. doi:10.1002/gdj3.20

1950

1951 Nelder and Mead (1965) *Computer Journal*, Vol 7, pp 308-313

1952

1953 Newman M., Alexander, M.A., Ault, T.R., Cobb, K.M., Deser, C., Di Lorenzo, E., Mantua, N.J.,
1954 Miller, A.J., Minobe, S., Nakamura, H., Schneider, N., Vimont, D.J., Phillips, A.S., Scott, J.D.
1955 and Smith, C.A. (2016), The Pacific decadal oscillation, revisited, *Journal of Climate*, 29(12),
1956 4399-4427, <https://doi.org/10.1175/JCLI-D-15-0508.1>

1957

1958 O’Carroll, A.G., Eyre, J.R. and Saunders, R.W. (2008), Three-Way Error Analysis between
1959 AATSR, AMSR-E, and In Situ Sea Surface Temperature Observations. *J. Atmos. Oceanic*
1960 *Technol.*, 25, 1197–1207, <https://doi.org/10.1175/2007JTECHO542.1>

1961

1962 Povey, A. C. and Grainger, R. G. (2015) Known and unknown unknowns: uncertainty estimation
1963 in satellite remote sensing, *Atmos. Meas. Tech.*, 8, 4699-4718, [https://doi.org/10.5194/amt-8-](https://doi.org/10.5194/amt-8-4699-2015)
1964 [4699-2015](https://doi.org/10.5194/amt-8-4699-2015).

1965

1966 Rasmussen, C.E. and Williams, C.K.I. (2006) Gaussian Processes for Machine Learning, MIT
1967 Press, 2006. ISBN-10 0-262-18253-X, ISBN-13 978-0-262-18253-9.
1968
1969 Rayner, N. A., Parker, D. E., Horton, E. B., Folland, C. K., Alexander, L. V., Rowell, D. P.,
1970 Kent, E. C. and Kaplan, A. (2003), Global analyses of sea surface temperature, sea ice, and night
1971 marine air temperature since the late nineteenth century, *J. Geophys. Res.*, 108, 4407, doi:
1972 10.1029/2002JD002670, D14.
1973
1974 Rayner N.A., Brohan, P., Parker, D.E., Folland, C.K., Kennedy, J.J., Vanicek, M., Ansell, T.J.,
1975 and Tett, S.F.B. (2006), Improved analyses of changes and uncertainties in sea surface
1976 temperature measured in situ since the mid-nineteenth century: the HadSST2 data set, *Journal of*
1977 *Climate*, 19(3), 446-469, <https://doi.org/10.1175/JCLI3637.1>
1978
1979 Research Data Archive/Computational and Information Systems Laboratory/National Center for
1980 Atmospheric Research/University Corporation for Atmospheric Research, Physical Sciences
1981 Division/Earth System Research Laboratory/OAR/NOAA/U.S. Department of Commerce,
1982 Cooperative Institute for Research in Environmental Sciences/University of Colorado, National
1983 Oceanography Centre/University of Southampton, Met Office/Ministry of Defence/United
1984 Kingdom, Deutscher Wetterdienst (German Meteorological Service)/Germany, Department of
1985 Atmospheric Science/University of Washington, Center for Ocean-Atmospheric Prediction
1986 Studies/Florida State University, and National Centers for Environmental
1987 Information/NESDIS/NOAA/U.S. Department of Commerce. 2016, updated monthly.
1988 International Comprehensive Ocean-Atmosphere Data Set (ICOADS) Release 3, Individual

1989 Observations. Research Data Archive at the National Center for Atmospheric Research,
1990 Computational and Information Systems Laboratory. <https://doi.org/10.5065/D6ZS2TR3>.
1991 Accessed 12 August 2016.
1992
1993 Smith, T.M. and Reynolds, R.W. (2002), Bias corrections for historical sea surface temperatures
1994 based on marine air temperatures, *Journal of Climate* 15(1), 73-87, [https://doi.org/10.1175/1520-
1995 0442\(2002\)015<0073:BCFHSS>2.0.CO;2](https://doi.org/10.1175/1520-0442(2002)015<0073:BCFHSS>2.0.CO;2)
1996
1997 Smith, T.M., Reynolds, R.W., Peterson, T.C., and Lawrimore, J. (2008), Improvements to
1998 NOAA's Historical Merged Land–Ocean Surface Temperature Analysis (1880–2006). *J.*
1999 *Climate*, 21, 2283–2296, <https://doi.org/10.1175/2007JCLI2100.1>
2000
2001 Sybrandy A.L., Niiler, P.P., Martin, C., Scuba, W., Charpentier, E. and Meldrum, D.T. (2008),
2002 Global Drifter Programme Barometer Drifter Design Reference, DBCP Report No. 4
2003 REVISION 2.2, August 2009. Published by the *Data Buoy Cooperation Panel*. Downloaded
2004 from <http://www.jcommops.org/dbcp/community/standards.html> 17 February 2018.
2005
2006 Thompson, D.W.J., Kennedy, J.J., Wallace, J.M. and Jones, P.D. (2008), A large discontinuity in
2007 mid 20th century global-mean surface temperatures, *Nature* 453, 646-649,
2008 <http://dx.doi.org/10.1038/nature06982>
2009

2010 Thorne, P.W., Parker, D.E., Christy, J.R. and Mears, C.A. (2005), Uncertainties in climate trends
2011 - lessons from upper- air temperature records, *Bulletin of the American Meteorological Society*
2012 86(10), 1437-1442, <https://doi.org/10.1175/BAMS-86-10-1437>

2013

2014 Wilcox, R. (2001), *Fundamentals of Modern Statistical Methods: Substantially Improving Power*
2015 *and Accuracy*, New York: Springer.

2016

2017 Willett, K.M., Jones, P.D., Gillett, N.P., and Thorne, P.W. (2008), Recent Changes in Surface
2018 Humidity: Development of the HadCRUH Dataset. *J. Climate*, 21, 5364–5383,
2019 <https://doi.org/10.1175/2008JCLI2274.1>

2020

2021 Woodruff S.D., Worley, S.J., Lubker, S.J., Ji, Z., Freeman, J.E., Berry, D.I., Brohan, P., Kent,
2022 E.C., Reynolds, R.W., Smith, S.R. and Wilkinson, C. (2011), ICOADS release 2.5: extensions
2023 and enhancements to the surface marine meteorological archive, *International Journal of*
2024 *Climatology*, 31, 951-967, <https://doi.org/10.1002/joc.2103>

2025

2026 Xu, F. and Ignatov, A. (2010), Evaluation of in situ sea surface temperatures for use in the
2027 calibration and validation of satellite retrievals, *Journal of Geophysical Research: Oceans* (1978-
2028 2012), 115, C09022, doi: 10.1029/2010JC006129.

2029

2030 **Tables**

2031 Table 1: Values used for uncertainties arising from uncorrelated random effects (σ_u) and from
 2032 random effects which are correlated for a particular agent (σ_b) for ships, drifting buoys and
 2033 moored buoys from K11c.

	σ_u	σ_b
Ships	0.74°C	0.71°C
Drifting buoys	0.26°C	0.29°C
Moored buoys	0.30°C	0.20°C

2034

2035 **Figure Captions**

2036 Figure 1: numbers of observations passing QC (a, c, e) and super-observations (b, d, f, see
 2037 Section 3 for the definition) per month for the globe (black), Southern Hemisphere (orange) and
 2038 Northern Hemisphere (blue) for (a,b) 1850-1880, (c, d) 1880-2000 and (e, f) 2000-2018. Note
 2039 the very different scales for the y-axes.

2040

2041 Figure 2: Fractional contribution of different SST observation methods, 1915-2018, to (a) the
 2042 Global average, (b) the Southern Hemisphere average and (c) the Northern Hemisphere average.
 2043 The brown/orange/tan areas indicate ship observations as labeled in panel (a) and the blue areas
 2044 indicate buoy observations. The pale lilac area represents unknown measurement method
 2045 (assumed to be from ships). These are the initial assignments (Section 2.2.1) and are not the
 2046 assignments finally used to calculate the adjustments.

2047 Figure 3: Example fields from the gridding procedure for the 5° by 5° by pentad grid boxes: (a)
 2048 SST anomalies (°C) for June 2003 relative to the 1961-1990 average; (b) number of observations

2049 contributing to each grid-box average; (c) number of super observations in each grid box; (d)
 2050 fractional contribution to the grid-box average by ERI measurements; (e) fractional contribution
 2051 to the grid-box average by drifting buoys (drifters); (f) fractional contribution to the grid-box
 2052 average by moored buoys (moorings); (g) estimated uncertainty associated with uncorrelated
 2053 errors; (h) estimated uncertainty ($^{\circ}\text{C}$) associated with under sampling; and (i) estimated
 2054 uncertainty ($^{\circ}\text{C}$) associated with micro-bias errors.

2055 Figure 4: Schematic diagram showing the flow of information in the parameter and bias
 2056 estimation and the corresponding sections in the paper. Blue boxes are input data sets and
 2057 analyses. Pink boxes correspond to steps described in Section 2 and 3. Orange boxes correspond
 2058 to processing described in Section 4.

2059 Figure 5: Annual global averages of the 200 realizations of the bucket corrections (a) R06-style
 2060 corrections (b) SR02-style corrections, (c) combined R06 and SR02 corrections. Blue lines in
 2061 (a)-(c) show ensemble members with a linear transition from wooden to canvas buckets and
 2062 orange lines show ensemble members generated assuming a step change in the fraction of
 2063 wooden and canvas buckets. (d) Estimated seasonal cycle of insulated bucket biases.

2064 Figure 6: (left column) Estimated seasonal-average ERI biases (orange) and bucket biases (blue)
 2065 ($^{\circ}\text{C}$) 1940-2018 for (a) the globe, (c) the southern hemisphere and (e) the northern hemisphere.
 2066 (right column) Smoothed estimated monthly ERI biases ($^{\circ}\text{C}$) for (b) the globe, (d) the southern
 2067 hemisphere and (f) the northern hemisphere. The smoothed estimate is shown in orange with the
 2068 seasonal-averages from the left column shown in grey.

2069 Figure 7: Smoothed time series of the estimated fraction f_{correct} of measurements labeled as
 2070 buckets that were correctly identified as buckets for (a) rejected start and end dates and (b)

2071 accepted start and end dates. (c) accepted (blue) and rejected (red) start and end date
2072 combinations. (d) The black line is the annual average inferred fraction of correct bucket
2073 assignments which is the average of the unsmoothed blue lines in (b). The red lines indicate the
2074 mean and ranges used to draw samples. The mean values are held constant at 0.5 between 1945
2075 and 1952 and at 0.95 after 1978. The blue lines show 10 samples of the full 200 member
2076 ensemble from 1945 on. Before 1945 the mean is set to 0.5, but the uncertainty is larger.

2077

2078 Figure 8: (a) Estimated monthly bias ($^{\circ}\text{C}$, orange) in global average SST 1850-2018 for the full
2079 gridded dataset, including ships and buoys and (b) estimated bias in global average SST anomaly
2080 ($^{\circ}\text{C}$, blue) relative to a 1961-1990 period, again for all data. (c and d) as for (a and b) except for
2081 the Southern Hemisphere. (e and f) as for (a and b) except for the Northern Hemisphere. The
2082 black line is the median bias and the shaded area represents the 95% range of the estimated
2083 biases.

2084 Figure 9: Monthly global average SST anomalies ($^{\circ}\text{C}$) 1850-2018 (a) relative to the unadjusted
2085 1961-1990 climatology and (b) relative to the bias-adjusted 1961-1990 climatology. (c and d) as
2086 for (a and b) except for the Southern Hemisphere. (e and f) as for (a and b) except for the
2087 Northern Hemisphere. The grey line shows the unadjusted data, the black line is the median of
2088 the adjusted data. The blue and orange shading represents the 95% range of the ensemble.

2089 Figure 10: (a) Annual global average unadjusted SST anomalies 1940-2014 ($^{\circ}\text{C}$ relative to
2090 unadjusted 1961-1990 climatology) for collocated bucket (blue) and ERI (orange)
2091 measurements. The solid line is the best estimate and the shaded area is the 95% uncertainty
2092 range (accounting for measurement and sampling errors). (b) Adjusted anomalies with expanded

2093 uncertainty range including bias adjustment uncertainty. The dotted line indicates the best
2094 estimate of the unadjusted series from (a). (c) and (d) as for (a) and (b) but for the Southern
2095 Hemisphere. (e) and (f) as for (a) and (b) but for the Northern Hemisphere.

2096 Figure 11: Comparisons between sea-surface temperature data sets from different sources 1990-
2097 July 2018, different comparison data sets cover different periods. (a) Global average SST
2098 anomalies ($^{\circ}\text{C}$), relative to the 1961-1990 HadSST2 climatology, from Argo floats (purple and
2099 purple shading) and HadSST.4.0.0.0 (black solid line and grey shading) each reduced to their
2100 common coverage. Shading indicates the 95% uncertainty range. The unadjusted SST data are
2101 shown as a black dotted line. (b) Indicator of the number of 5° grid boxes in HadSST.4.0.0.0
2102 (pale grey) and in the Argo data set (purple) Data are plotted only for the overlap. (c) as for (a),
2103 but with ARC (in red, no uncertainty range shown) substituted for Argo. (d) as for (b) but with
2104 ARC substituted for Argo. (e) as for (a) but with buoys (in blue and blue shading) substituted for
2105 Argo. (f) as for (b) but with buoys substituted for Argo.

2106 Figure 12: Average SST difference ($^{\circ}\text{C}$) between HadSST.4.0.0.0 and the three instrumentally
2107 homogeneous data sets (a) buoys 1995-2018, (b) ARC 1995-2012, and (c) Argo 2000-2017. (d)
2108 Shows the difference between HadSST.4.0.0.0 and the unadjusted gridded SSTs, 1995-2018.

2109 Figure 13: (a) Collocated global annual average NMAT anomalies ($^{\circ}\text{C}$) 1900-2010 offset by
2110 0.15°C (blue, relative to 1961-1990) and global annual average SST anomalies from
2111 HadSST.4.0.0.0 (black is central estimate and grey shading indicates 95% uncertainty range). (b)
2112 Offset NMAT anomalies minus SST anomalies with combined 95% uncertainty range (taking
2113 into account the bias errors from the HadSST.4.0.0.0 ensemble, and measurement and sampling
2114 errors in the SST). The dashed line indicates zero difference. (c) Collocated global annual

2115 average offset NMAT anomalies (blue), global annual average near-surface water temperature
2116 from HadIOD excluding Argo (orange) and SST (black is central estimate and grey shading
2117 again indicates 95% uncertainty range). (d) Difference between HadIOD and HadSST.4.0.0.0
2118 (orange) and NMAT and HadSST.4.0.0.0 (blue). The shaded area indicates the 95% uncertainty
2119 range.

2120 Figure 14: (a and b) Global, (c and d) Southern Hemisphere and (e and f) Northern Hemisphere
2121 annual average SST anomalies ($^{\circ}\text{C}$) 1850-2018 relative to 1961-1990 for HadSST.4.0.0.0 (black
2122 line is the median in the left column and the grey shading in the right column represents the 95%
2123 uncertainty range) and HadSST.3.1.1.0 (blue line is the median and the blue shading represents
2124 the 95% uncertainty range). Uncertainty estimates combine the bias-adjustment uncertainties
2125 from the ensemble with measurement and sampling uncertainties.

2126 Figure 15: (a) Global average SST anomaly 1850-2012 ($^{\circ}\text{C}$ relative to 1961-1990) series from
2127 HadSST.4.0.0.0 (black), ERSSTv5 (blue, thick line is operational version and thin-thick dashed
2128 lines are ensemble range from the 1000-member ERSSTv4 ensemble), COBE-SST-2 (orange),
2129 HadSST.3.1.1.0 (green) and unadjusted SSTs (red). All data sets are averaged on to a 5° grid and
2130 reduced to HadSST.4.0.0.0 coverage before comparison. (b) Global-average difference for each
2131 data set from HadSST.4.0.0.0. The grey shading shows the 95% uncertainty range from
2132 HadSST.4.0.0.0 including effects from measurement, sampling and bias-adjustment errors. The
2133 bias-adjustment uncertainty range is shown in darker grey. (c) and (d) as for (a) and (b) but for
2134 the Southern Hemisphere. (e) and (f) as for (a) and (b) but for the Northern Hemisphere.

2135 Figure 16: Global and regional average SST anomaly trends to 2012. Median trends from
2136 HadSST.4.0.0.0 are indicated by a black horizontal line and the grey shading indicates median

2137 and 95% uncertainty range including effects from measurement, sampling and bias-adjustment
2138 errors. The bias-adjustment uncertainty range is shown in darker grey; ERSSTv5 (blue), the
2139 lozenge is the operational version and the vertical line is the 95% ensemble range from the 1000-
2140 member ERSSTv4 ensemble); COBE-SST-2 (orange); HadSST.3.1.1.0 (green); and unadjusted
2141 SSTs (red).

2142 Figure 17: (a) Standard deviation of SST anomalies ($^{\circ}\text{C}$), (b) Zonal length scale L_x (km) and (c)
2143 meridional length scale L_y (km) used in the interpolation scheme..

2144

2145 Figure 18: Example of the inputs and outputs of the interpolation for July 2003. (a) gridded SST
2146 anomalies ($^{\circ}\text{C}$) from ships. (b) gridded SST anomalies from buoys. (c) interpolated SST
2147 anomalies. (d) estimated ERI biases ($^{\circ}\text{C}$). Note that the ERI data are a subset of the ship data so
2148 the coverage is not identical.

2149 Figure 19: Tests of effectiveness of the reconstruction method from 2000-2014. Blue lines show
2150 results for the 19th century coverage tests, red lines show the results of the missing-at-random
2151 tests and the purple lines show the test where data were removed at selected locations of Argo
2152 observations. (a) probability distribution for chi-squared statistics calculated from the withheld
2153 data. (b) average of histograms from all scaled, withheld observations in the 19th century
2154 coverage test (blue line) compared to the samples drawn from the posterior of the reconstruction
2155 (black line and grey shading). (c) as for (b) but for the missing-at-random test (red). (d) as for (b)
2156 but for the test where data were removed at the locations of Argo observations (purple). In (d)
2157 the dashed purple line shows the offset that occurs when doing the reconstruction based on

2158 unadjusted data and the thin blue line shows the scaled residuals of the adjusted SSTs at the
2159 locations of the Argo observations.

2160 Figure 20: Density map of estimated biases for ships from the HadSST.4.0.0.0 analysis (where
2161 the estimated uncertainty on the bias was less than 0.25°C) and from the IQUAM analysis
2162 between 2000 and 2012. The blue lines show the regression of the IQUAM estimate on
2163 HadSST.4.0.0.0 and vice versa. The red line is the $y=x$ line offset by 0.19°C . A small white cross
2164 marks $[0,0]$. White lines show the boundaries of the bins, separated by 0.1K .

2165 Figure 21: Tests of bias estimation using synthetic data. (a) histograms of differences between
2166 prescribed and estimated characteristic biases for each month for each measurement type in the
2167 synthetic data set including buckets (blue), ERI (orange) and hull sensors (hot pink). (b) time
2168 series of difference between assigned and estimated biases for bucket measurements (blue) and
2169 the 95% uncertainty range (grey lines). (c) as for (b) but for ERI measurements in orange. (d) as
2170 for (b) but for hull sensors in hot pink.

2171 Figure 22 (a) Density map of prescribed vs estimated micro-biases for all ships in the synthetic
2172 data set. Darker colours indicate higher densities. (b) as for (a) except it shows only those ships
2173 for which the uncertainty in the estimated micro-bias is less than 0.5°C . (c) as for (b) but with an
2174 uncertainty less than 0.15°C . (d) as for (b) but for uncertainties less than 0.05°C . The black
2175 diagonal line is $y=x$. White lines show the boundaries of the bins, separated by 0.1K .

2176 Figure 23: Histograms for trend estimates in global average temperatures ($^{\circ}\text{C}/\text{decade}$) for
2177 different periods (indicated in the individual plot titles) and different data sets. HadSST.4.0.0.0 is
2178 shown in grey (generated from the 200-member ensemble of bias adjustments only) and black
2179 (generated from the 200-member ensemble of bias adjustments combined with samples from the

2180 measurement and sampling uncertainties). The ERSSTv4 ensemble is shown in blue, with the
2181 operational ERSSTv5 point shown as a blue lozenge. COBE-SST-2 is shown in orange. The
2182 unadjusted data are shown in red. In addition, the trends from the IPCC AR5 SST data sets are
2183 shown as green lozenges numbered as follows: 1 is HadSST3, 2 is HadSST2, 3 is HadISST1.1, 4
2184 is COBE SST, and 5 is ERSSTv3. They do not appear in panels (e) and (f) as these periods were
2185 not considered in Chapter 2.

2186

Figure 1.

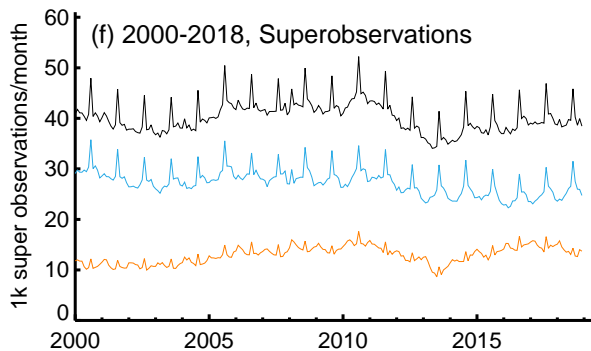
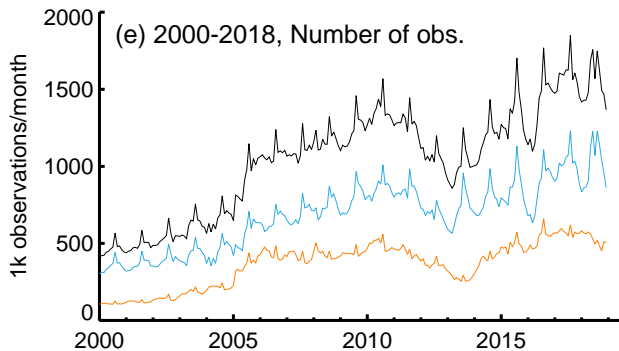
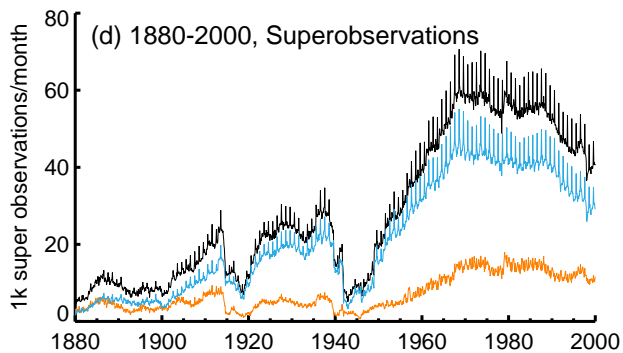
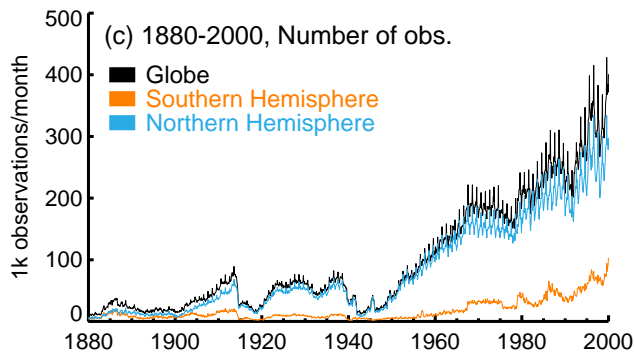
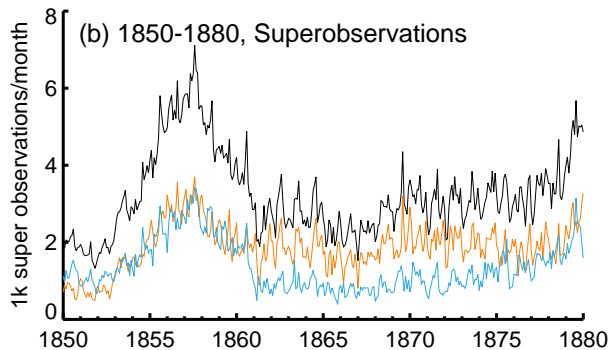
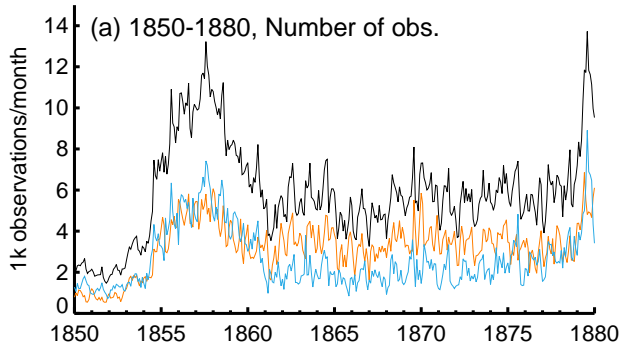


Figure 2.

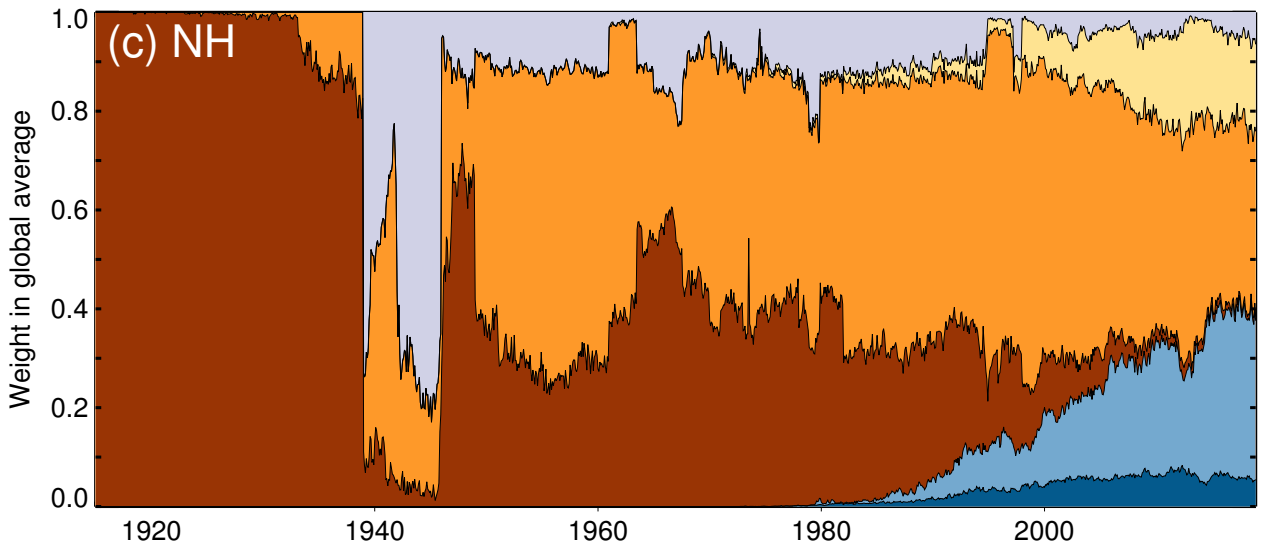
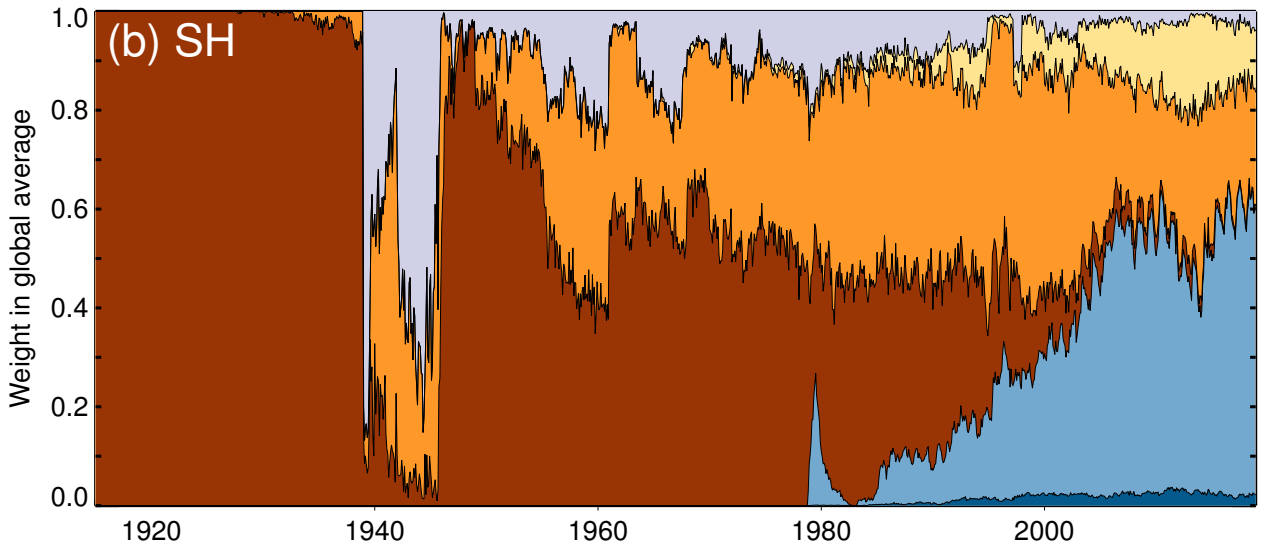
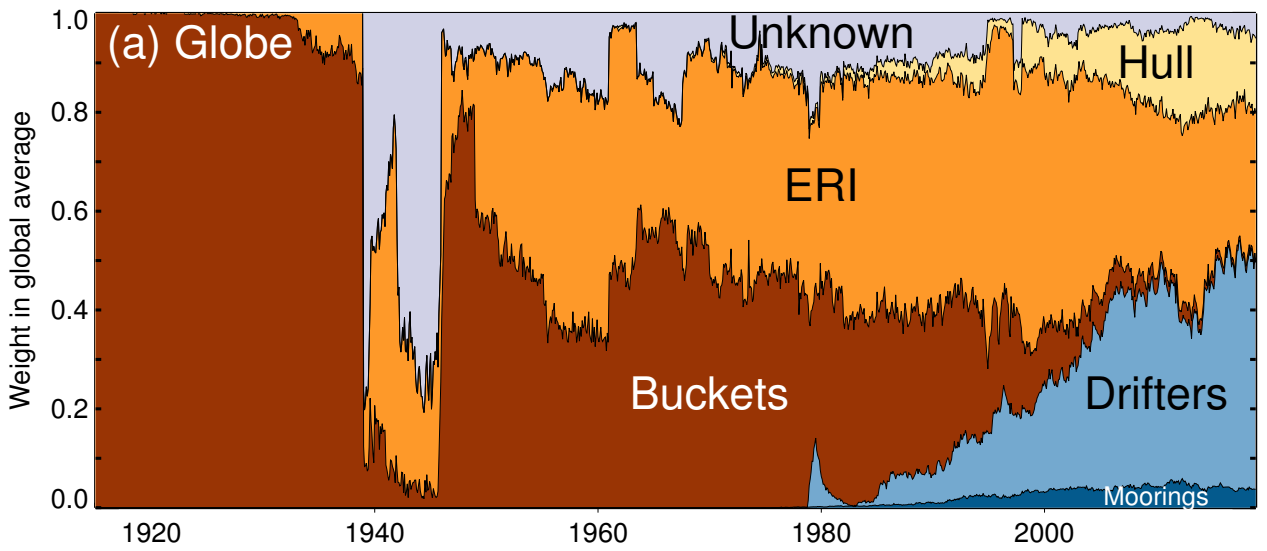
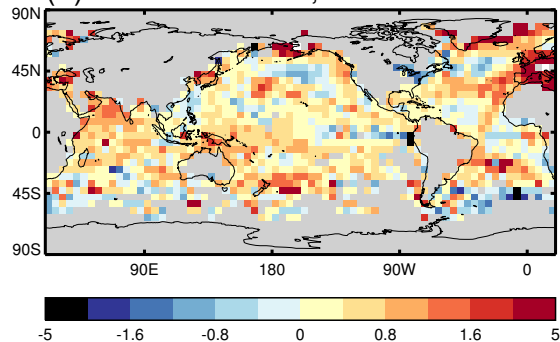
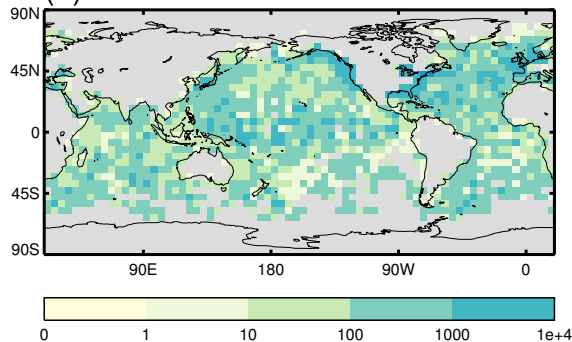


Figure 3.

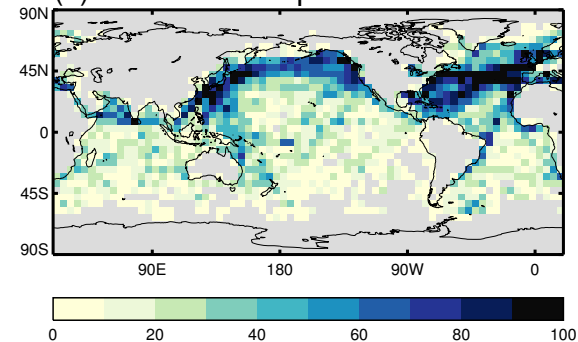
(a) SST anomalies, June 2003



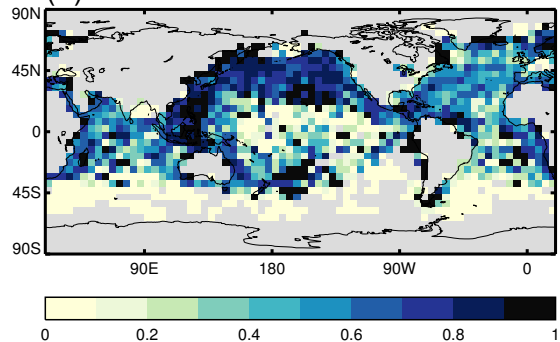
(b) Number of Observations



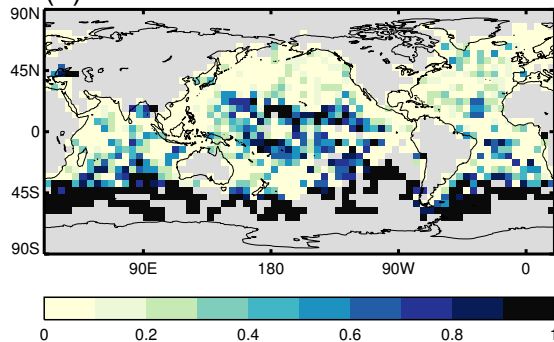
(c) Number of Super Observations



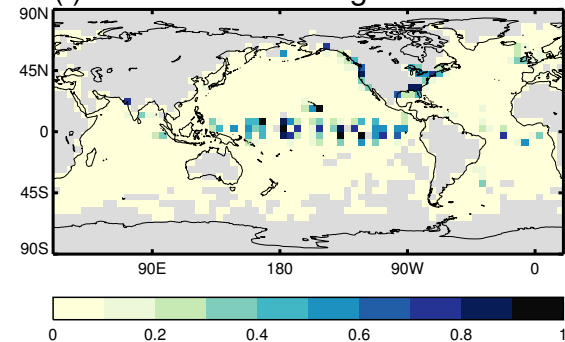
(d) Fraction of ERI



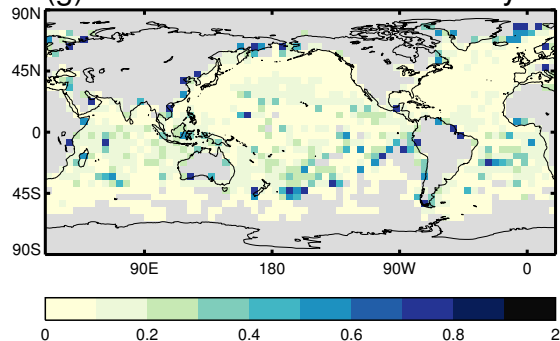
(e) Fraction of Drifters



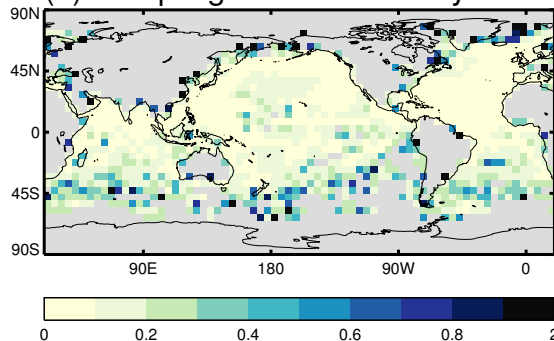
(f) Fraction of Moorings



(g) Uncorrelated error uncertainty



(h) Sampling error uncertainty



(i) Microbias error uncertainty

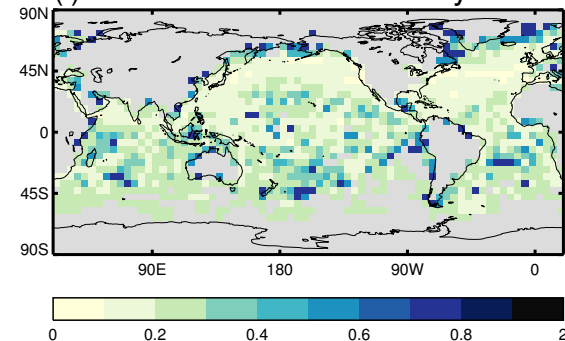


Figure 4.

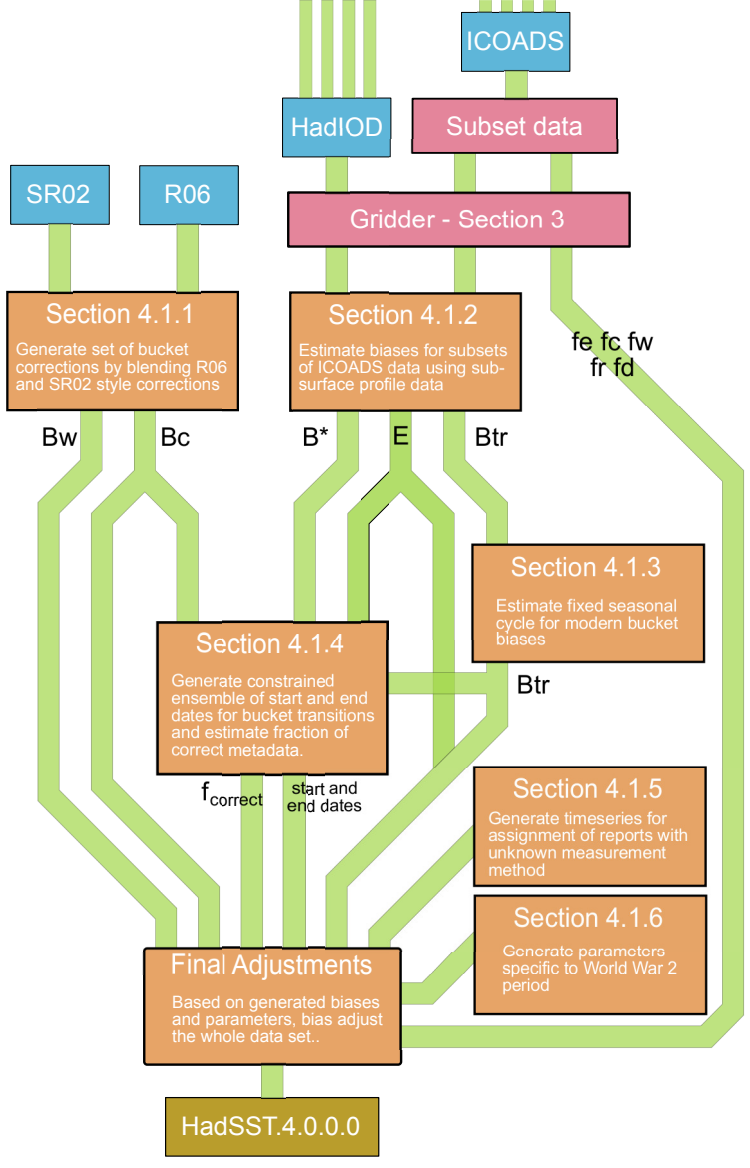


Figure 5.

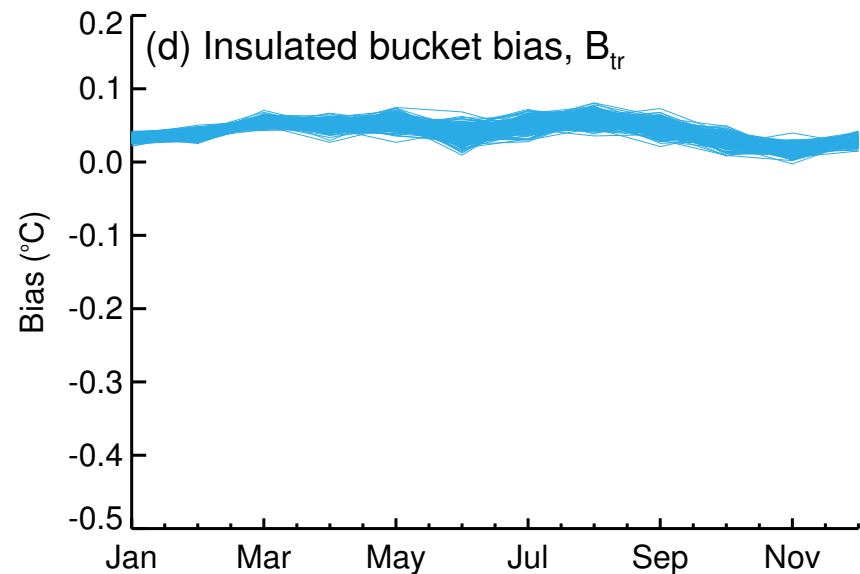
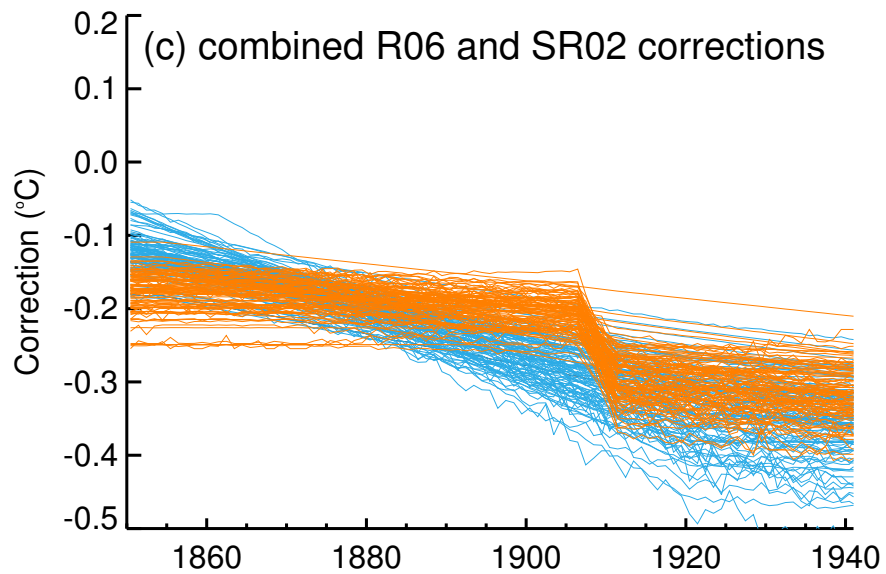
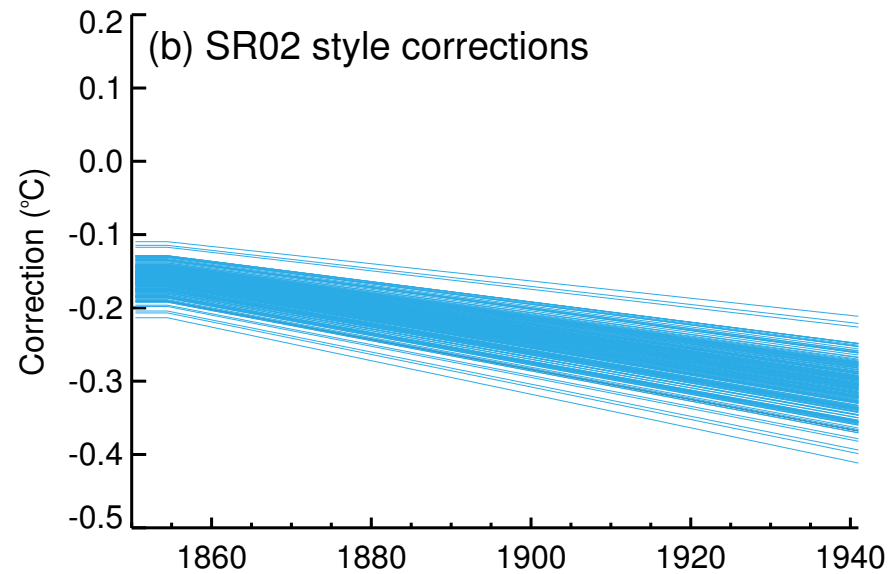
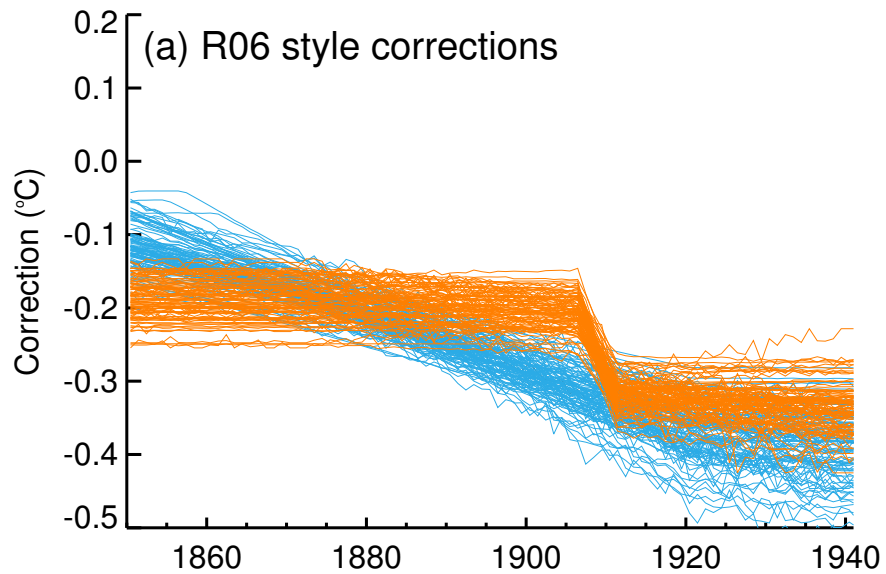


Figure 6.

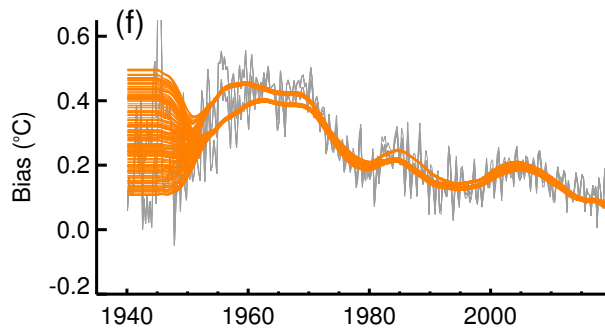
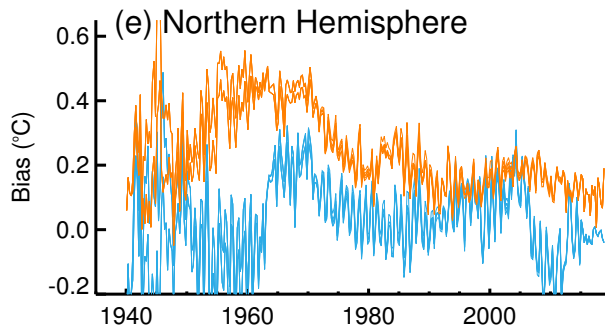
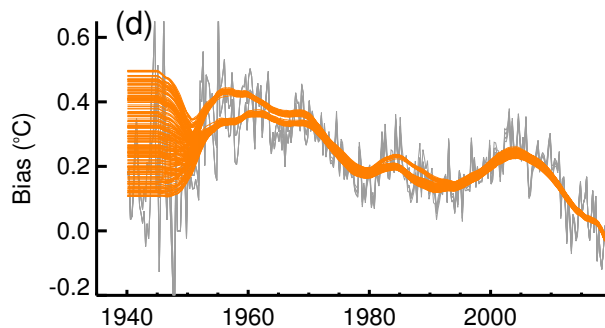
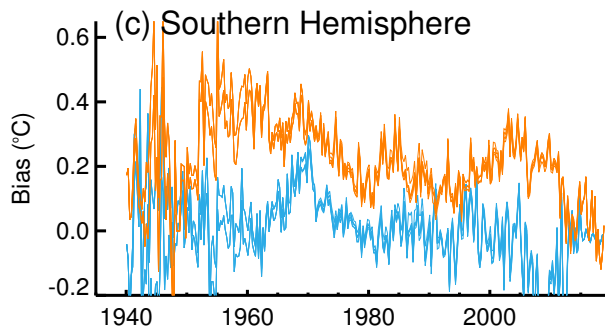
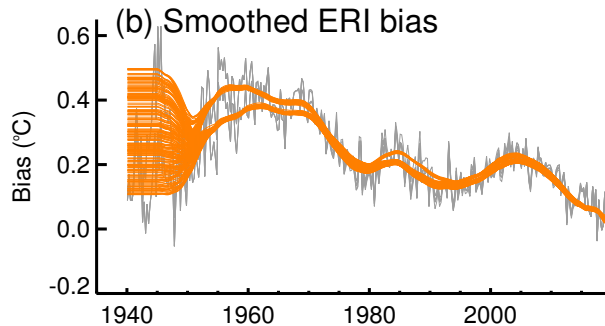
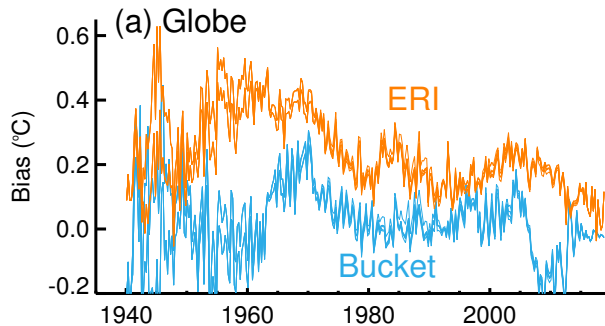


Figure 7.

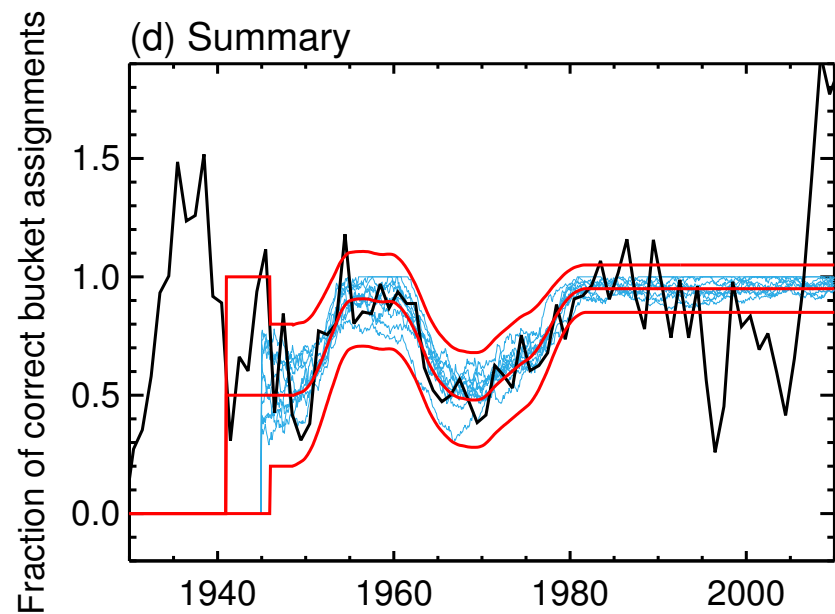
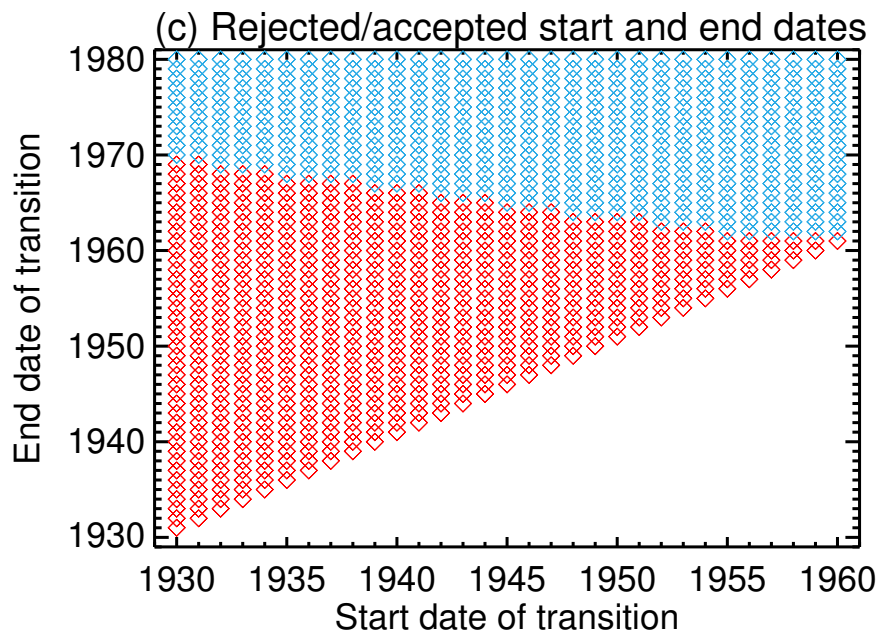
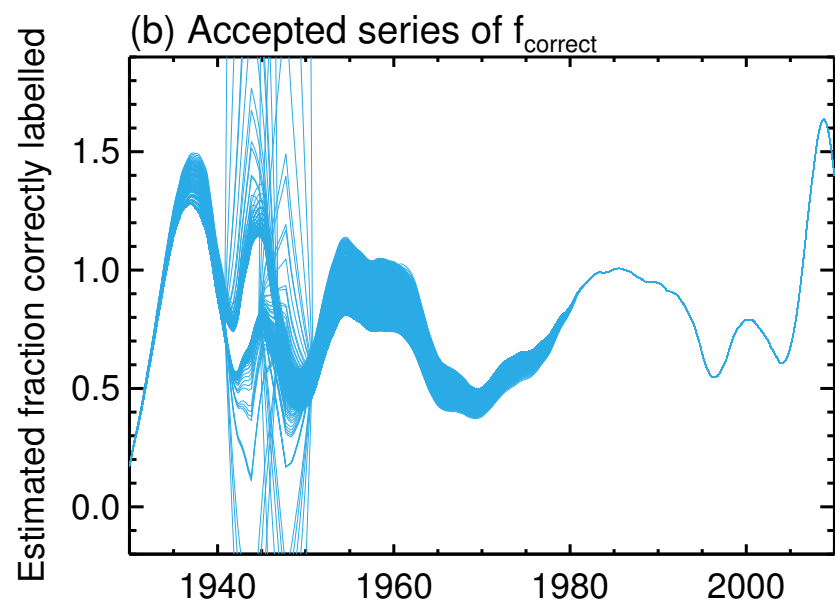
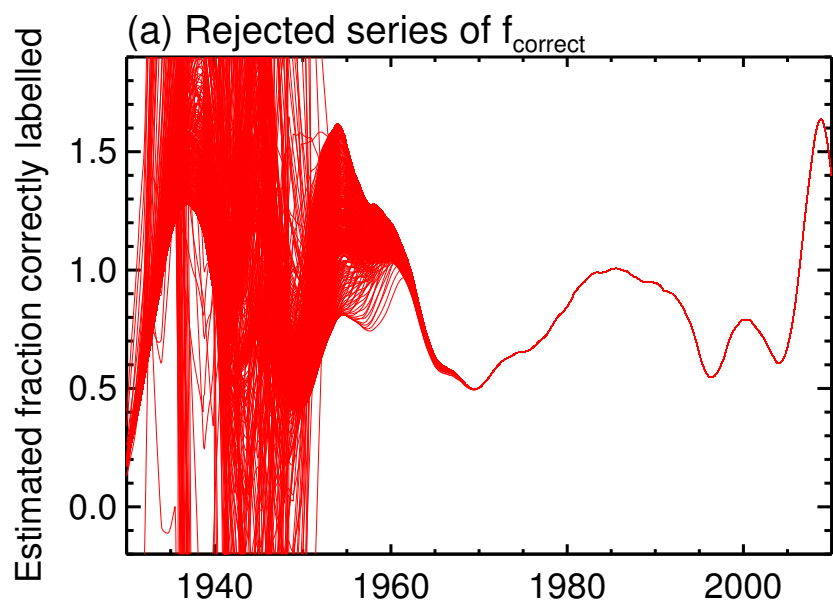


Figure 8.

Figure 9.

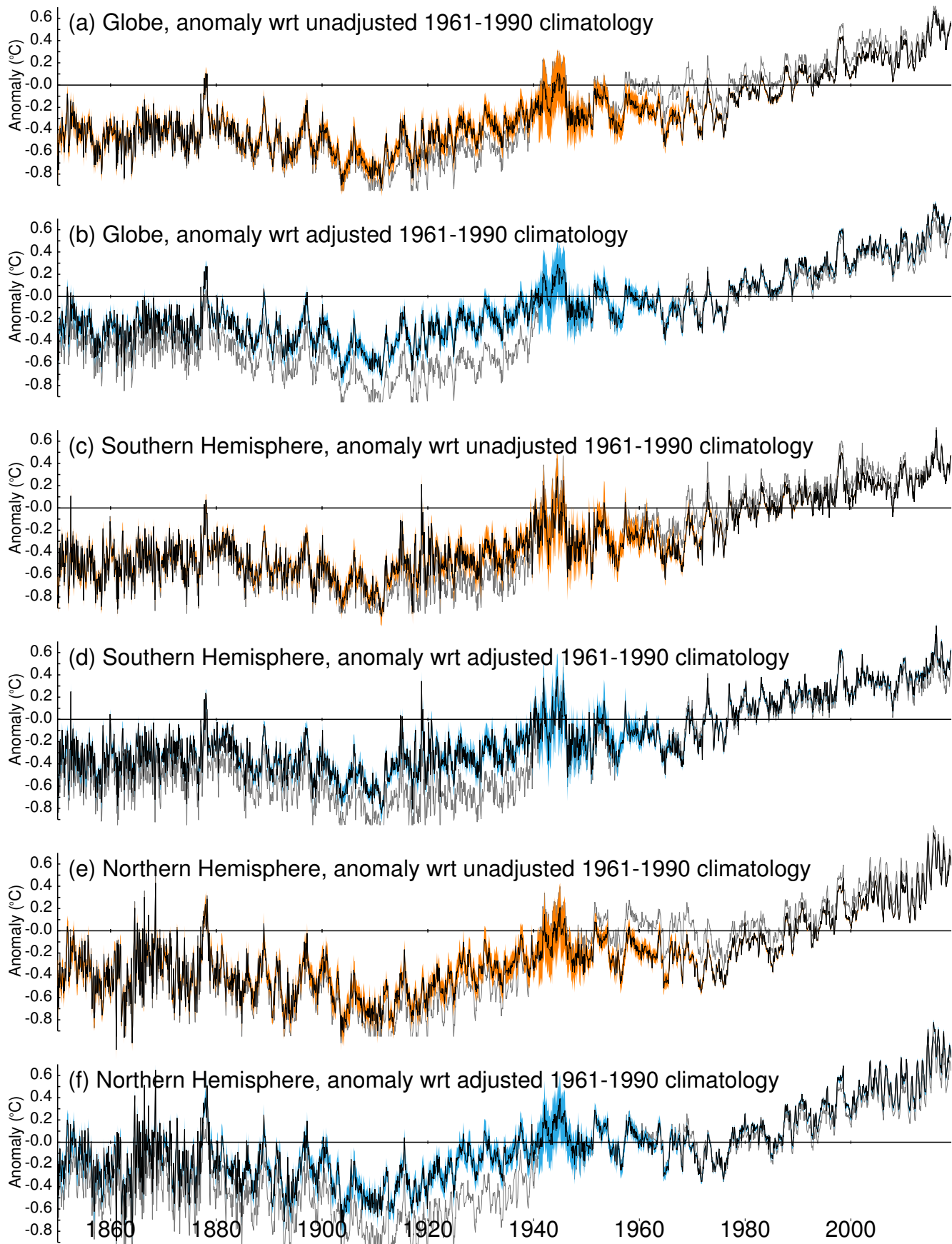


Figure 10.

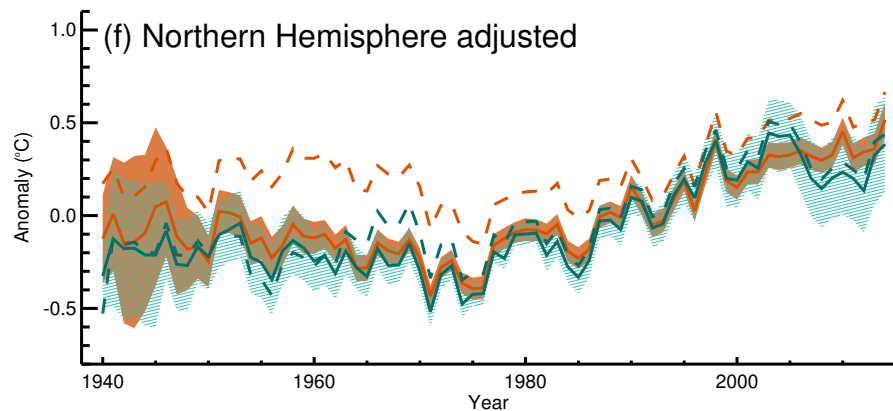
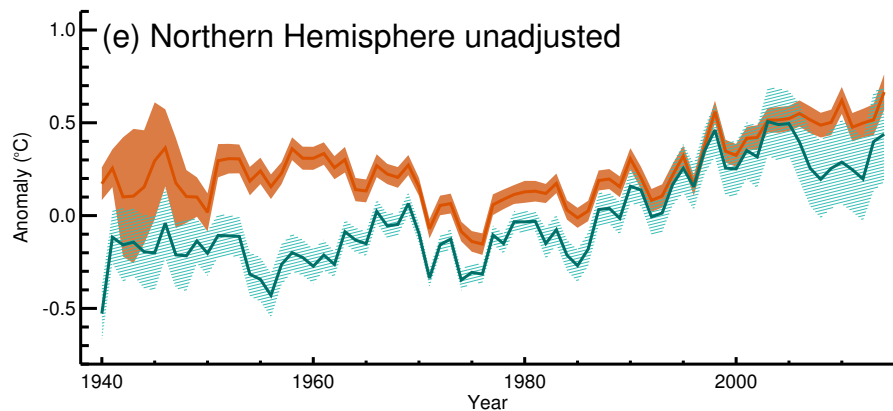
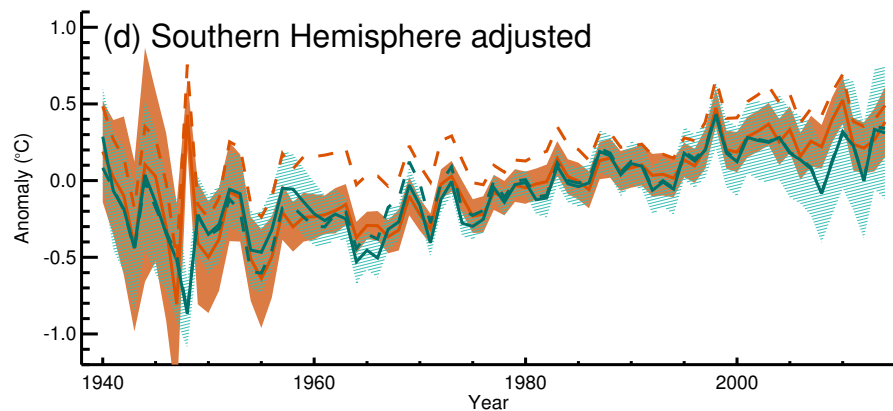
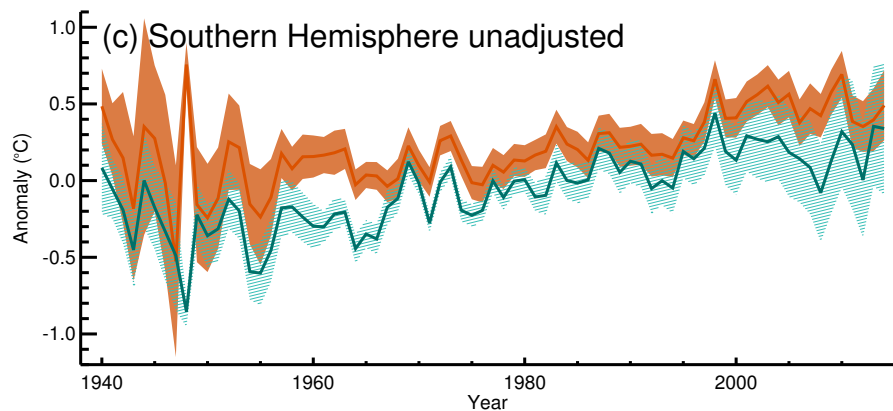
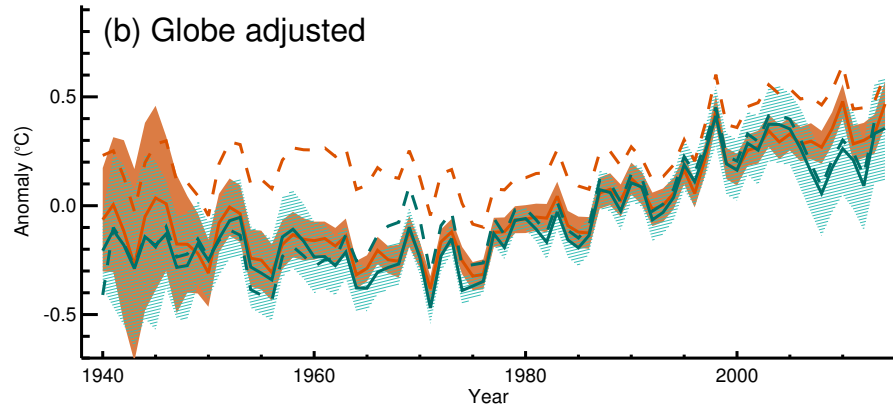
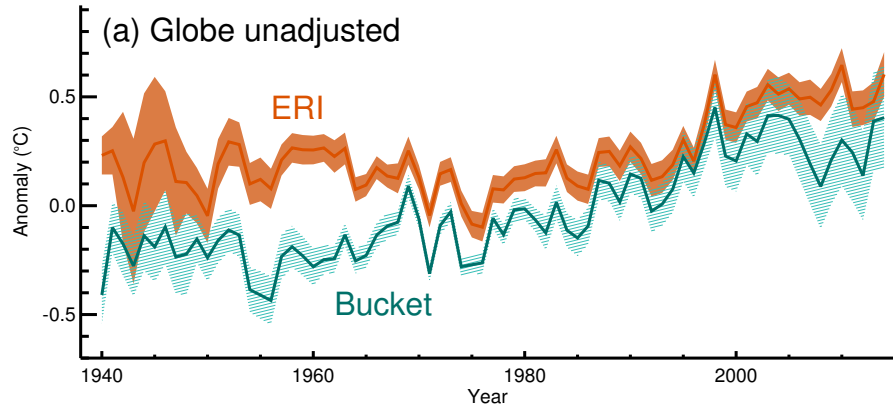


Figure 11.

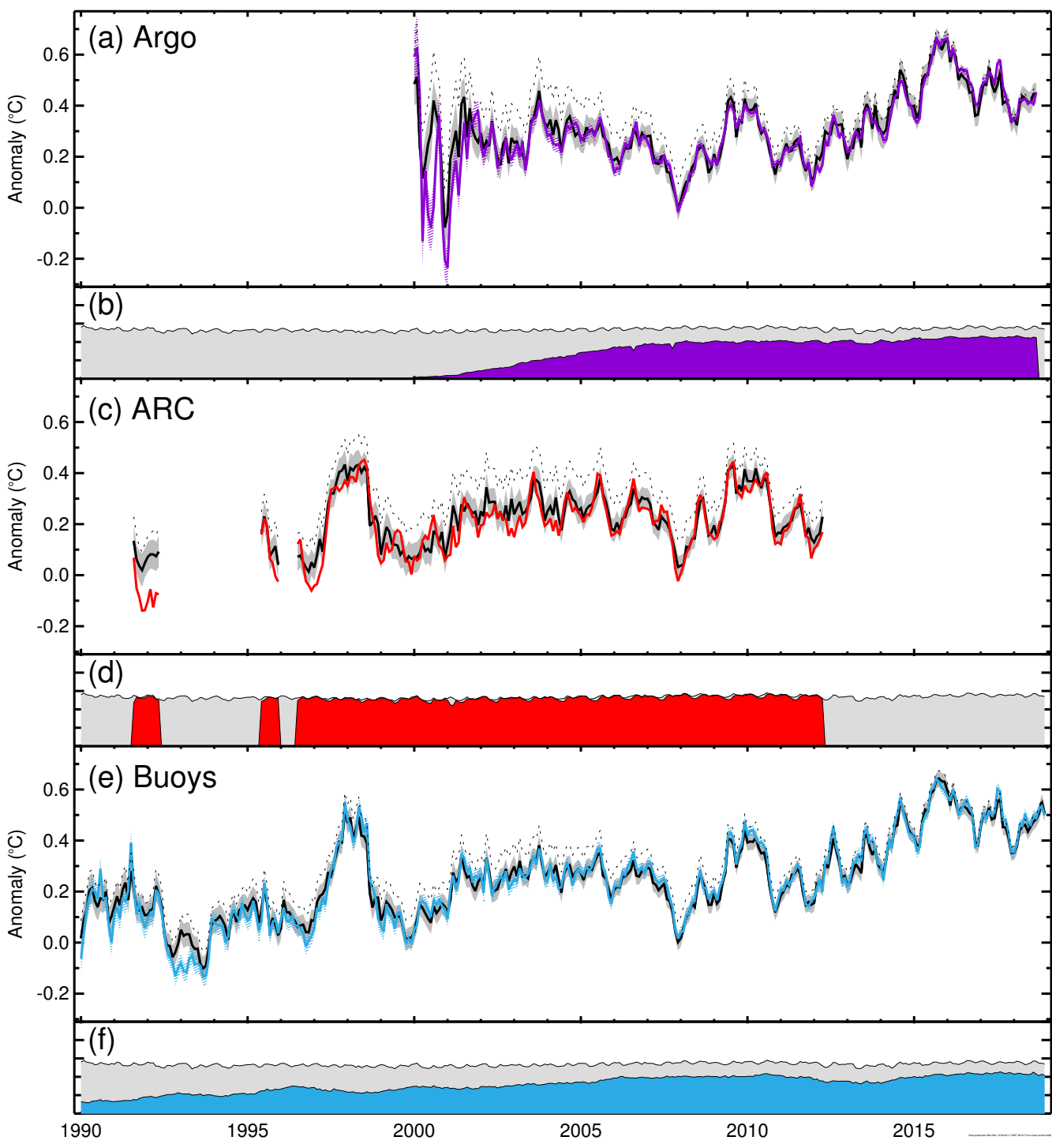
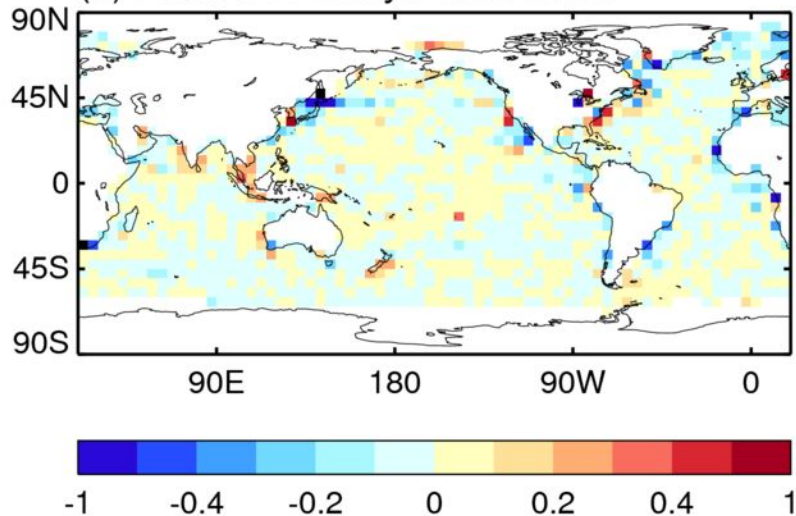
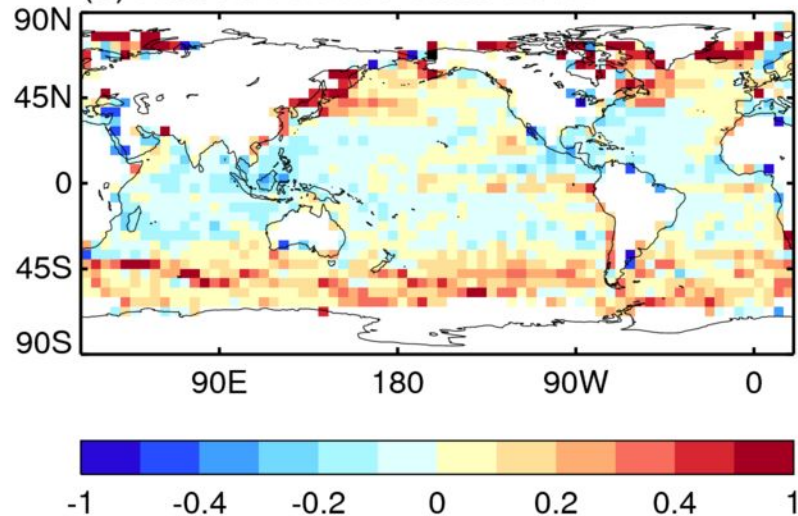


Figure 12.

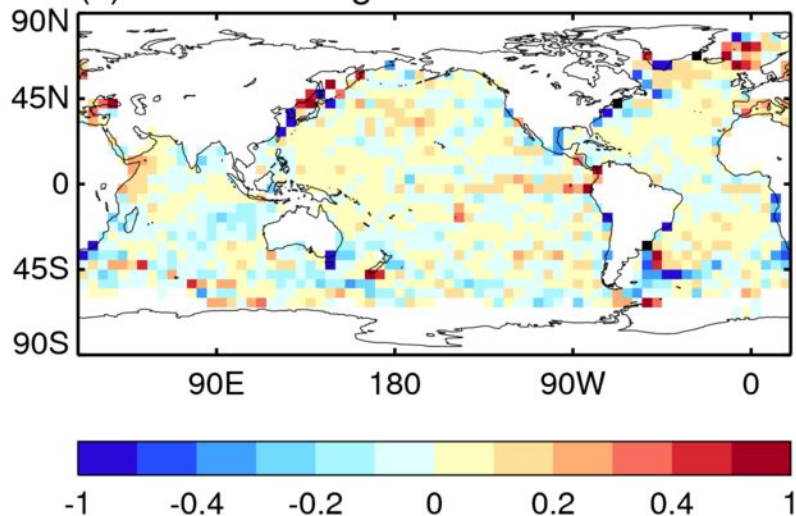
(a) HadSST4-Buoys 1995-2018



(b) HadSST4-ARC 1995-2012



(c) HadSST4-Argo 2000-2017



(d) HadSST4-Unadjusted data 1995-2018

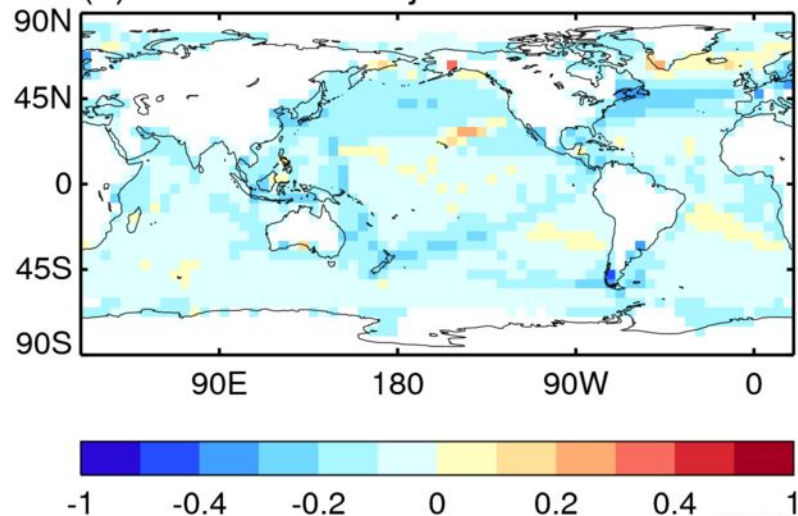


Figure 13.

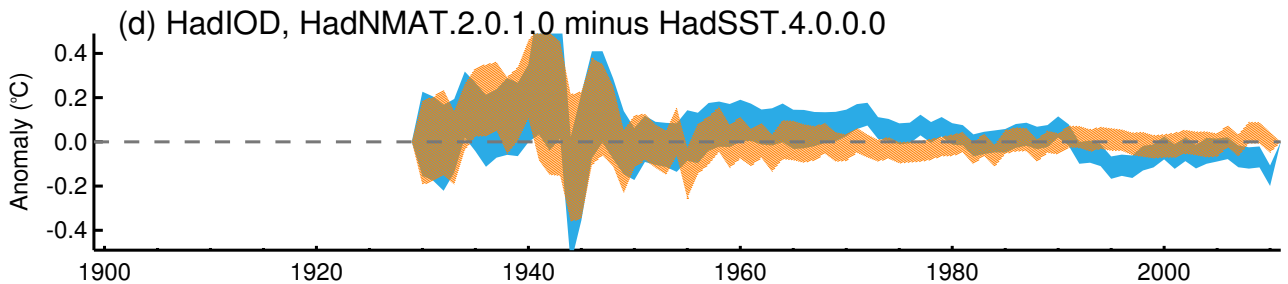
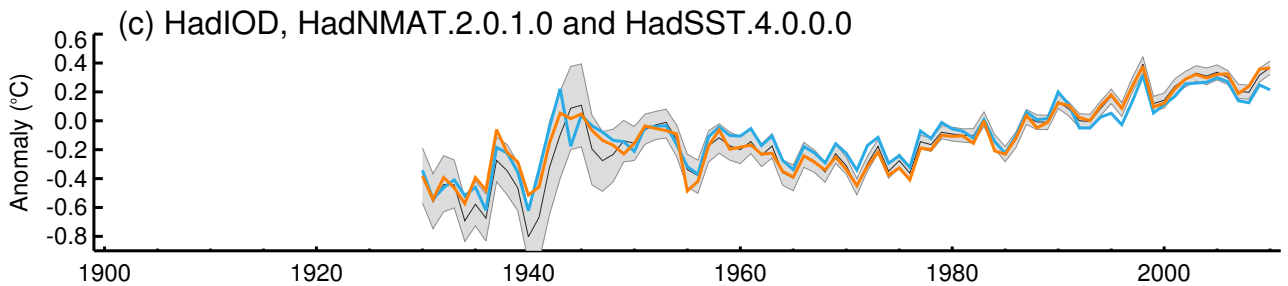
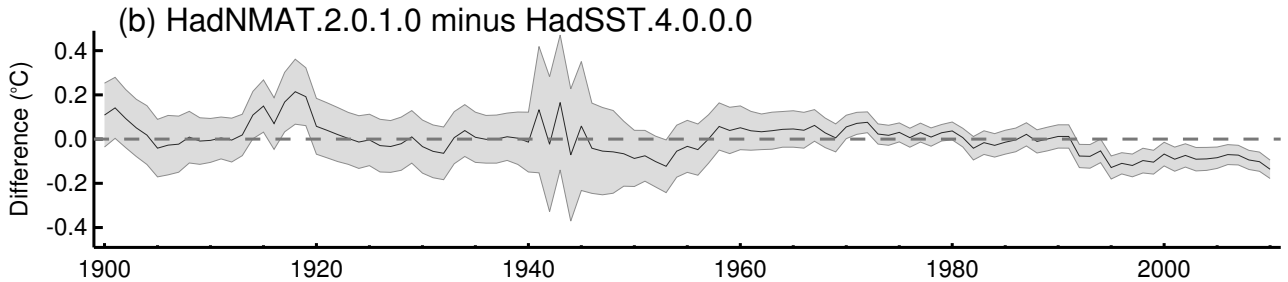
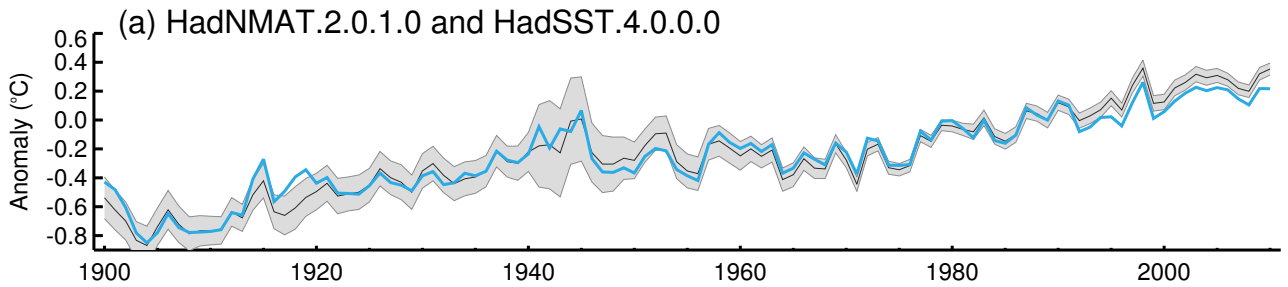


Figure 14.

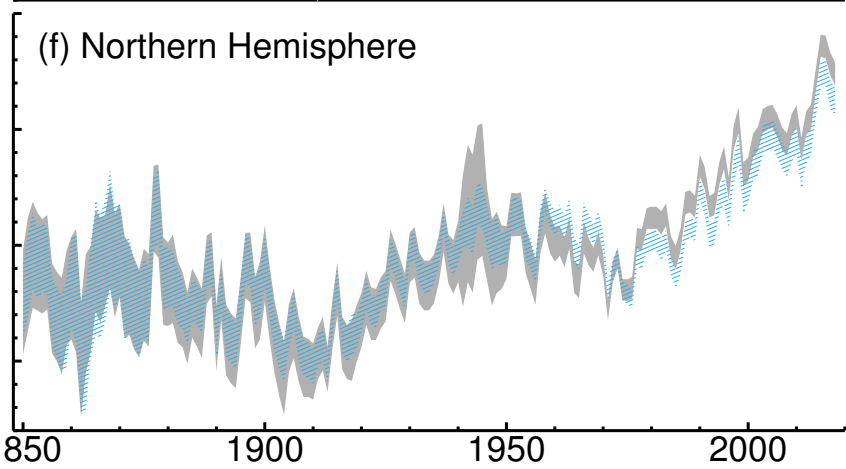
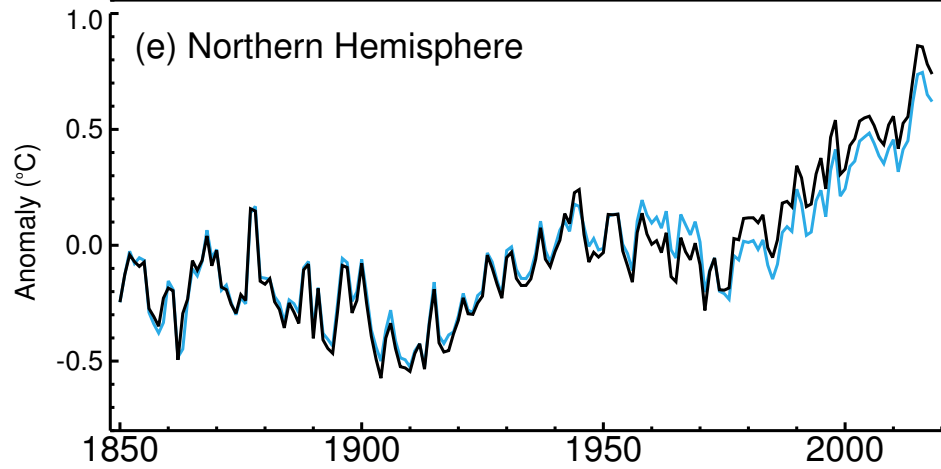
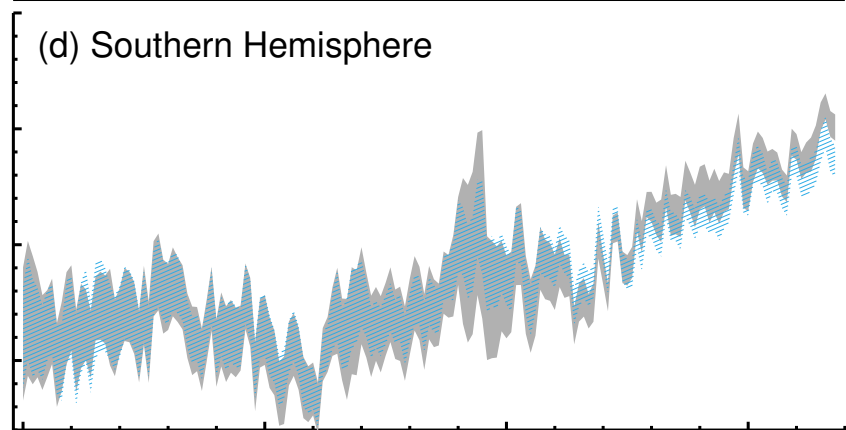
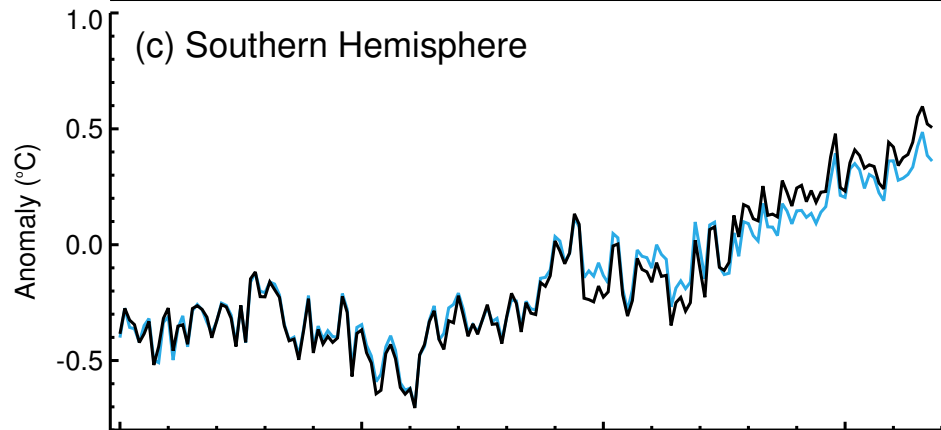
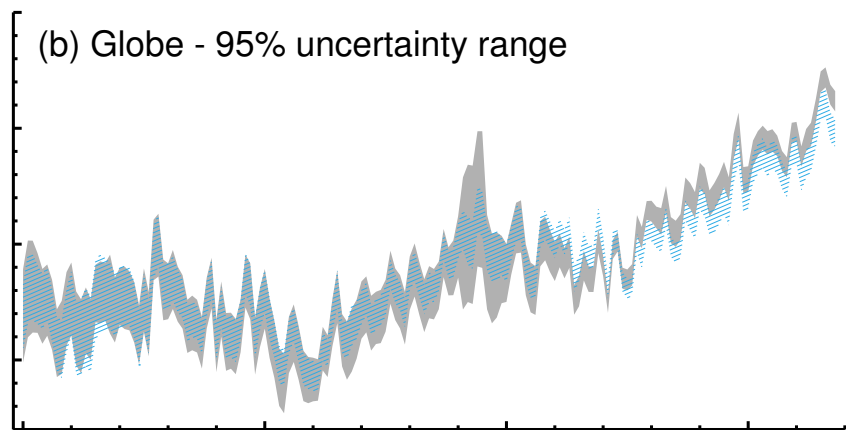
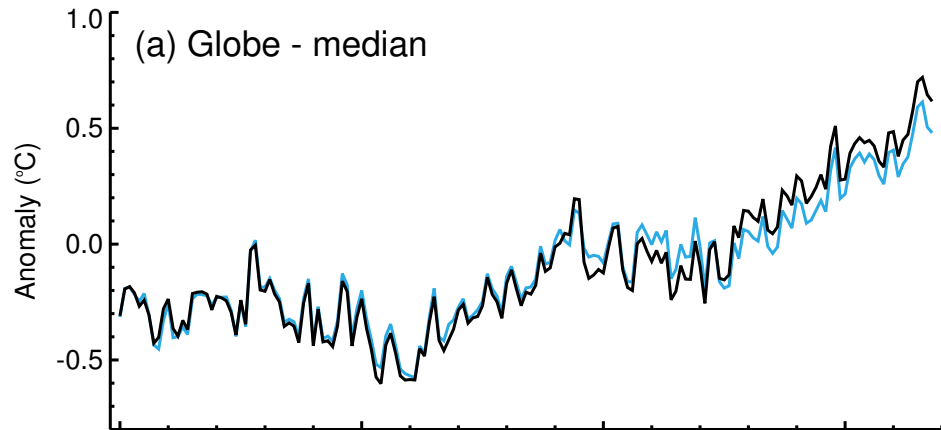


Figure 15.

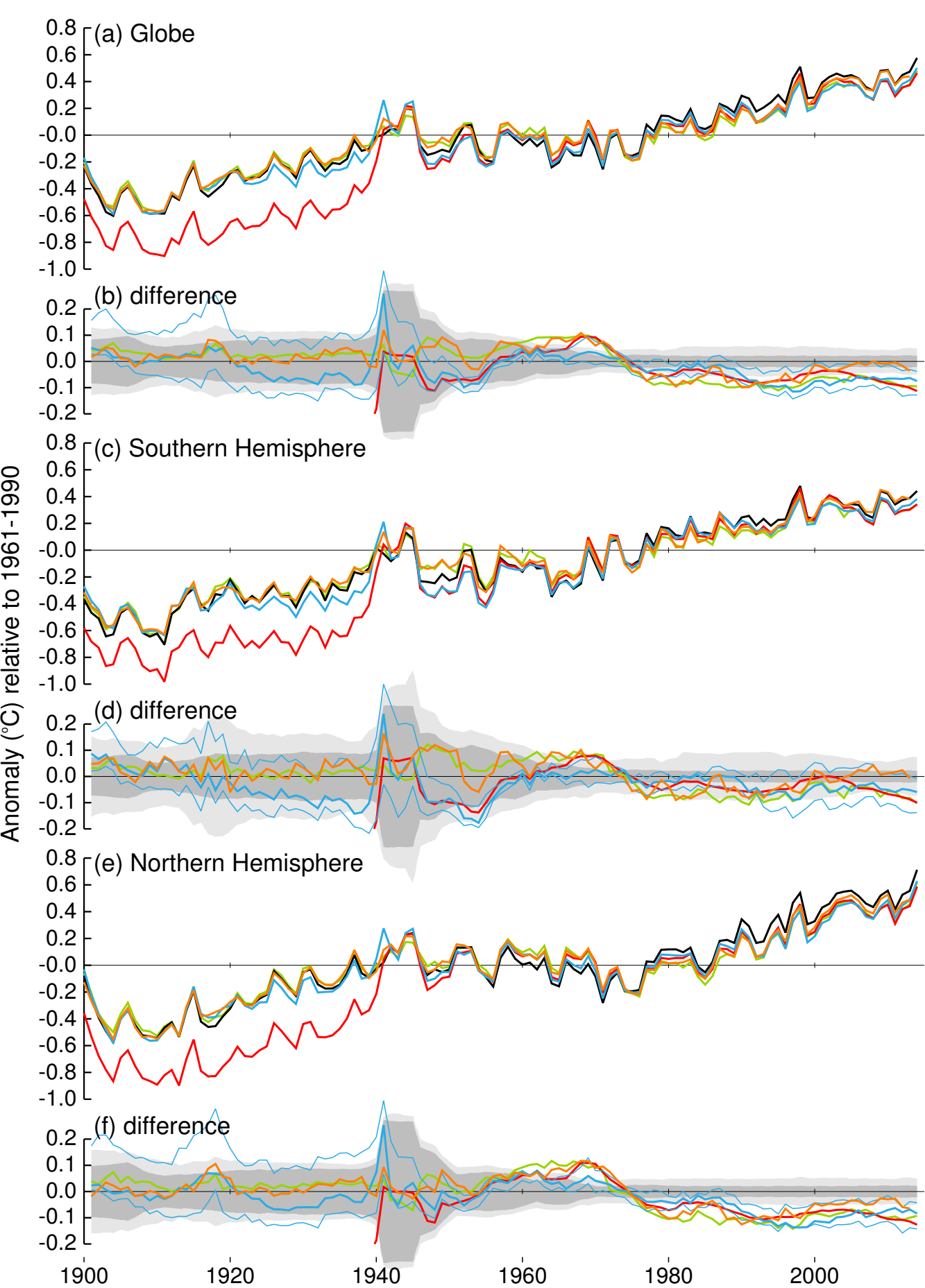


Figure 16.

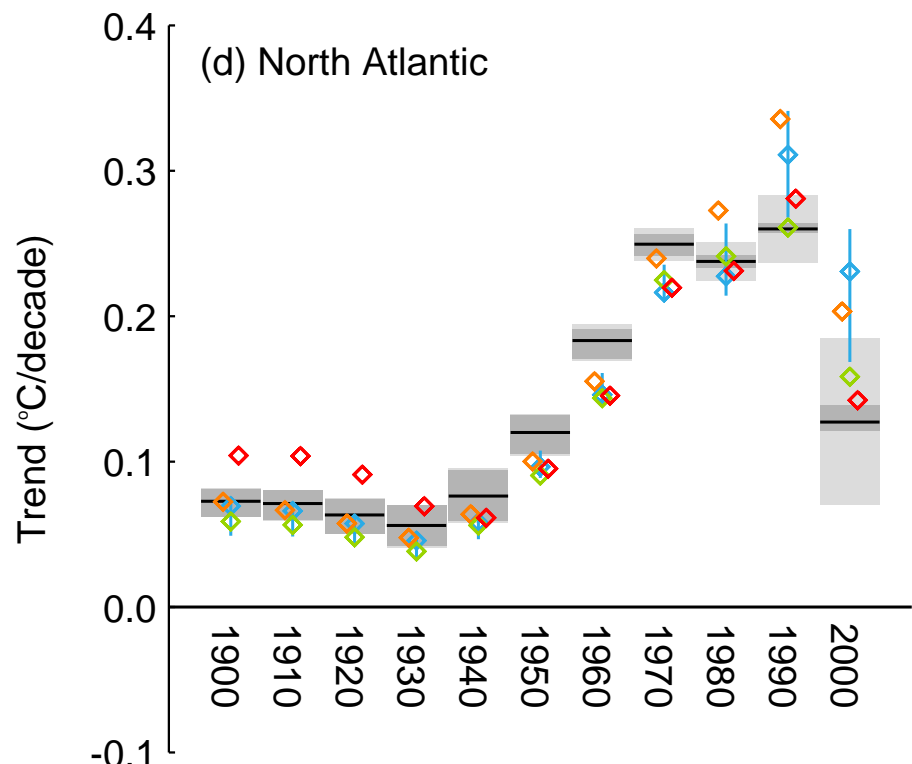
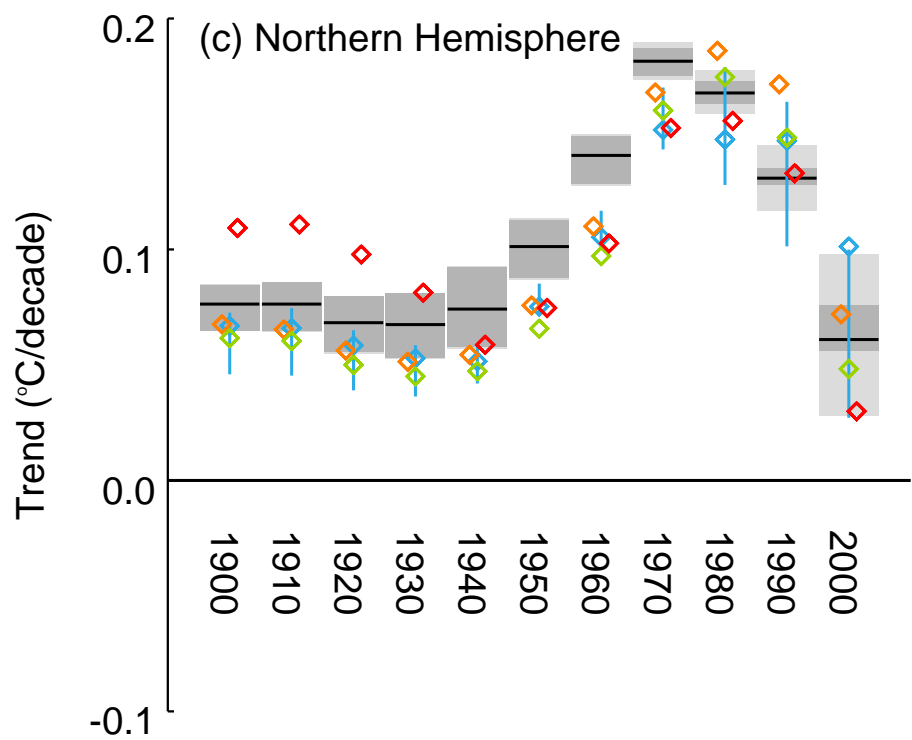
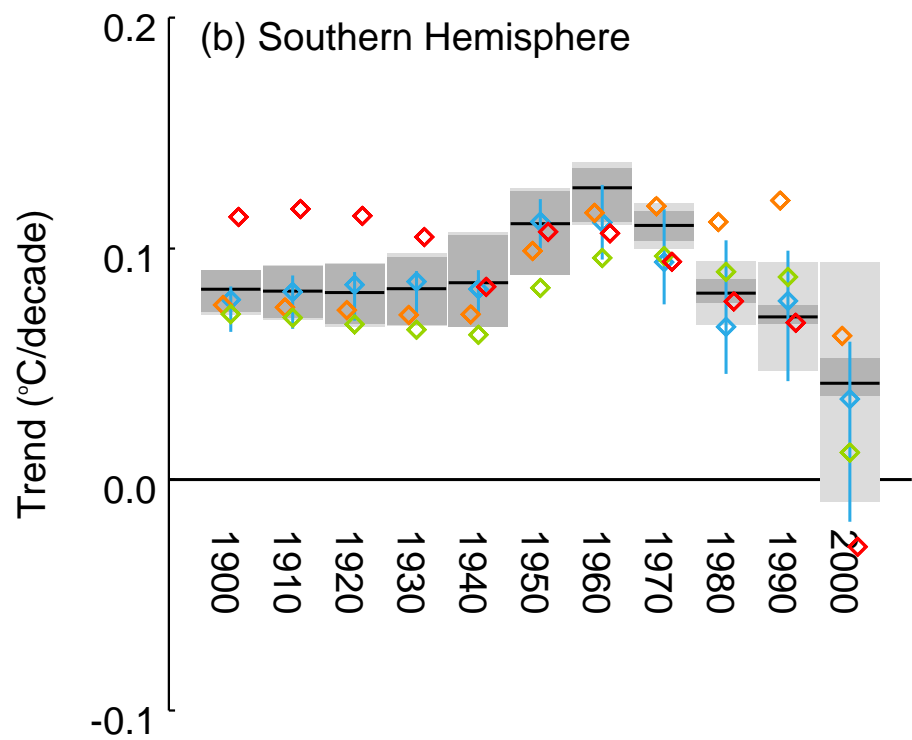
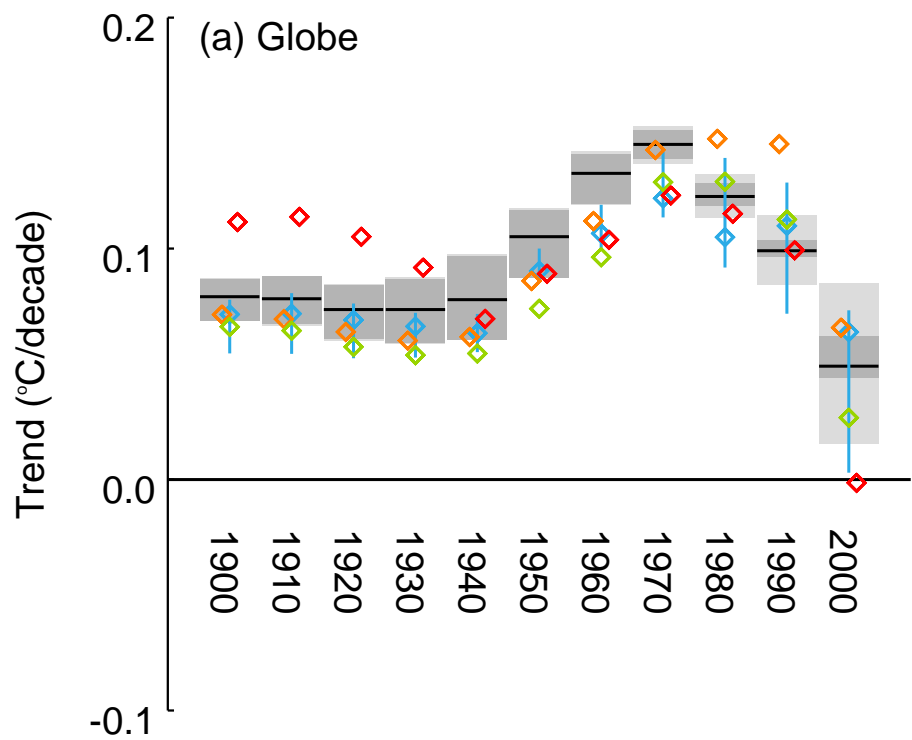
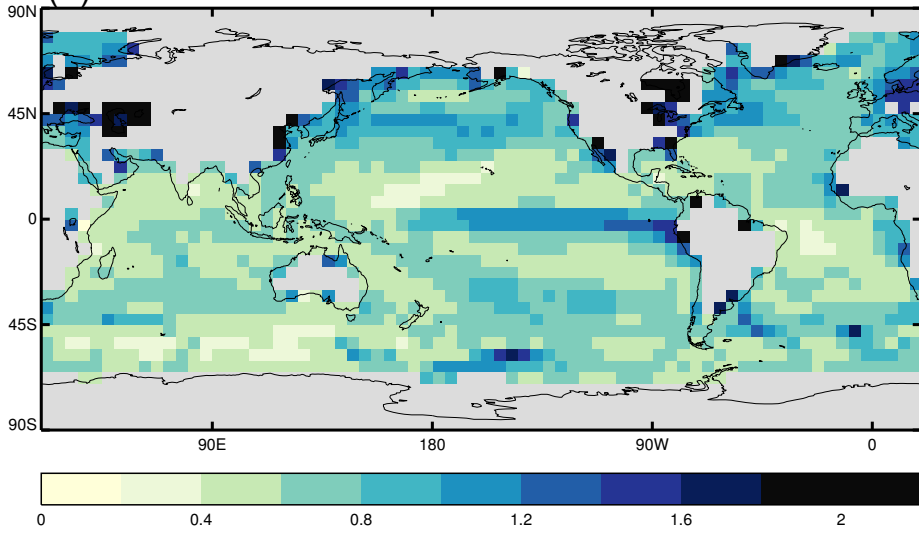
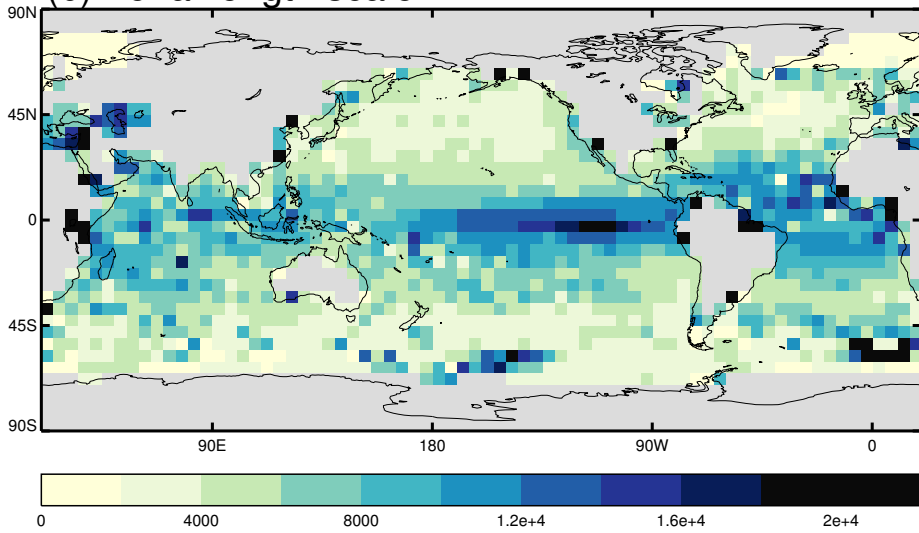


Figure 17.

(a) Standard Deviation



(b) Zonal length scale



(c) Meridional length scale

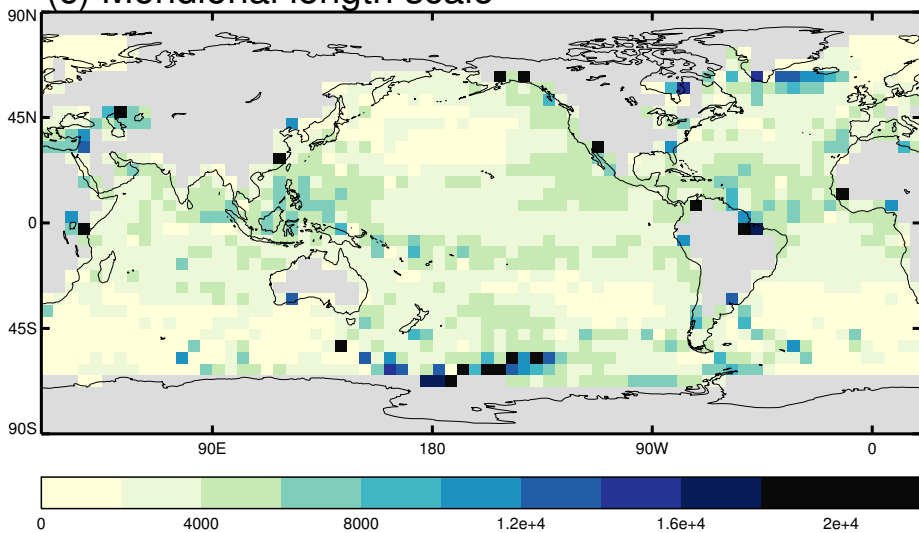
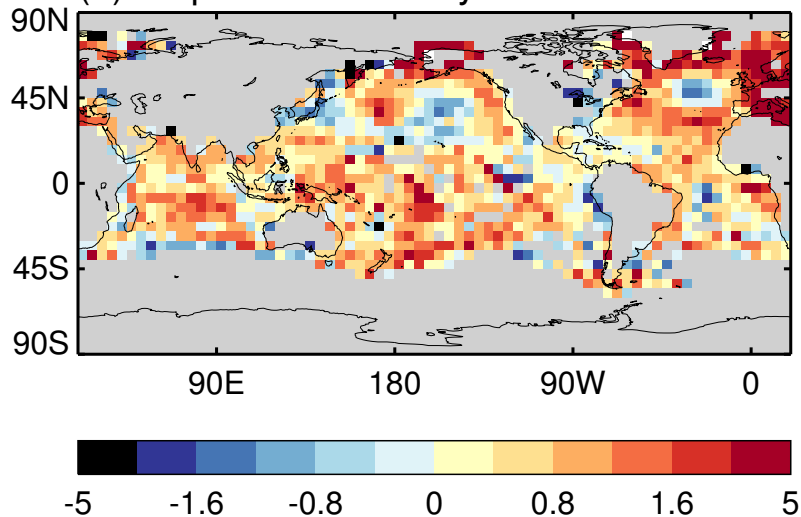
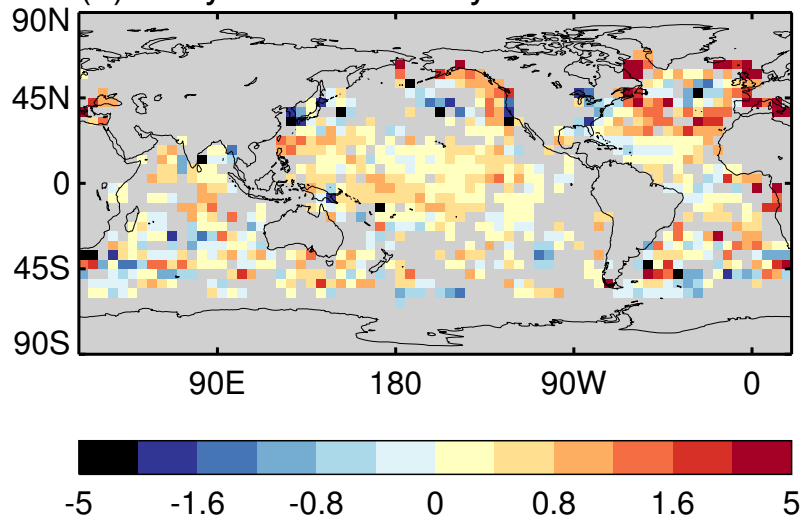


Figure 18.

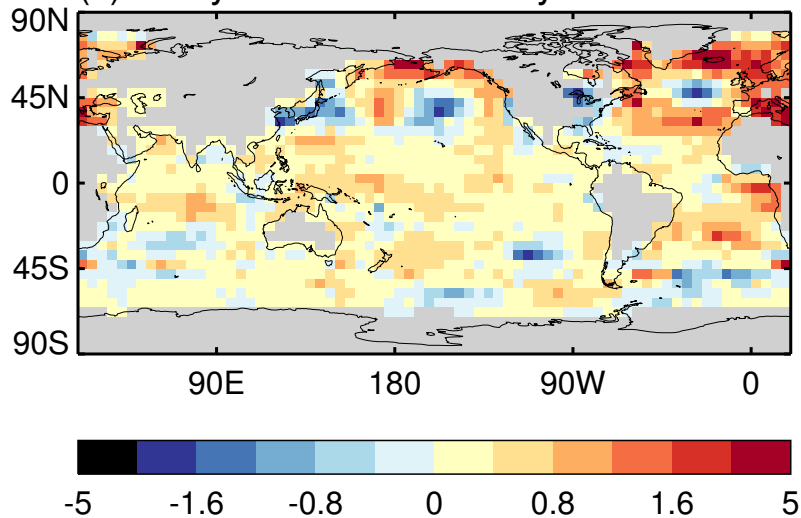
(a) Ship anomalies July 2003



(b) Buoy anomalies July 2003



(c) Analysed anomalies July 2003



(d) ERI Bias July 2003

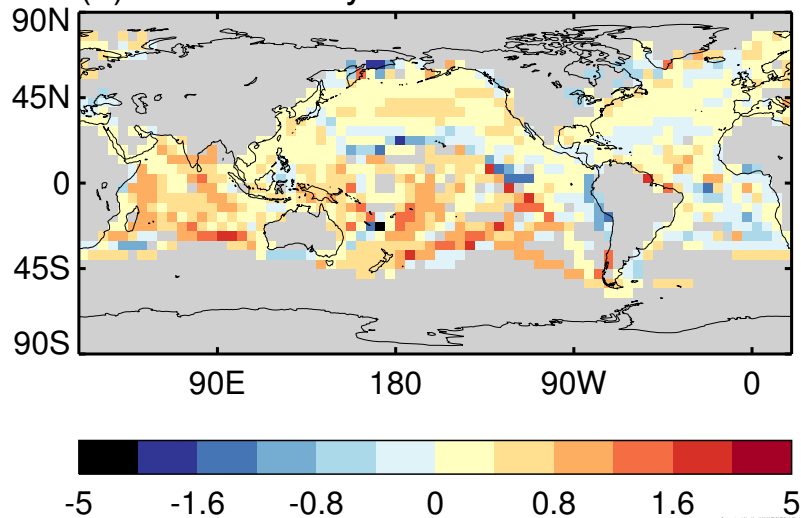


Figure 19.

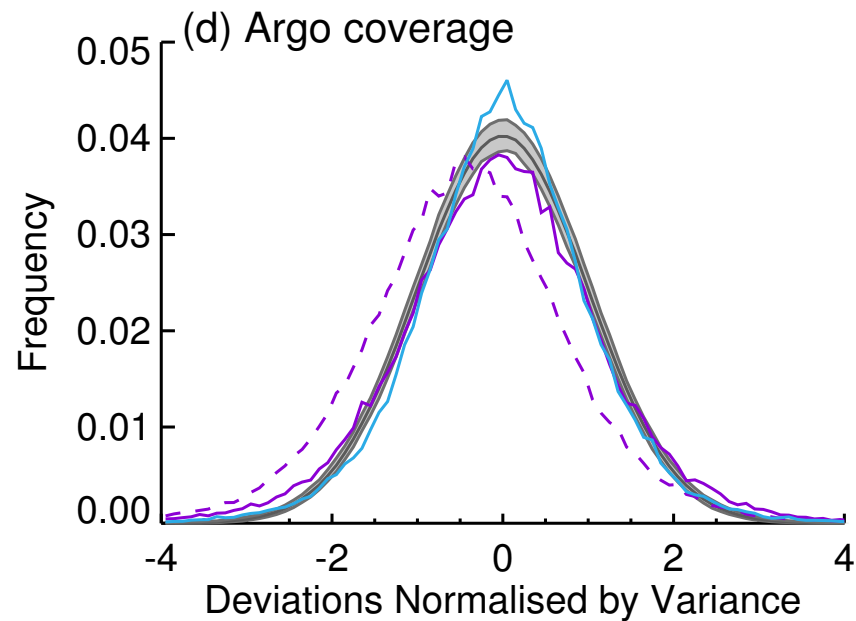
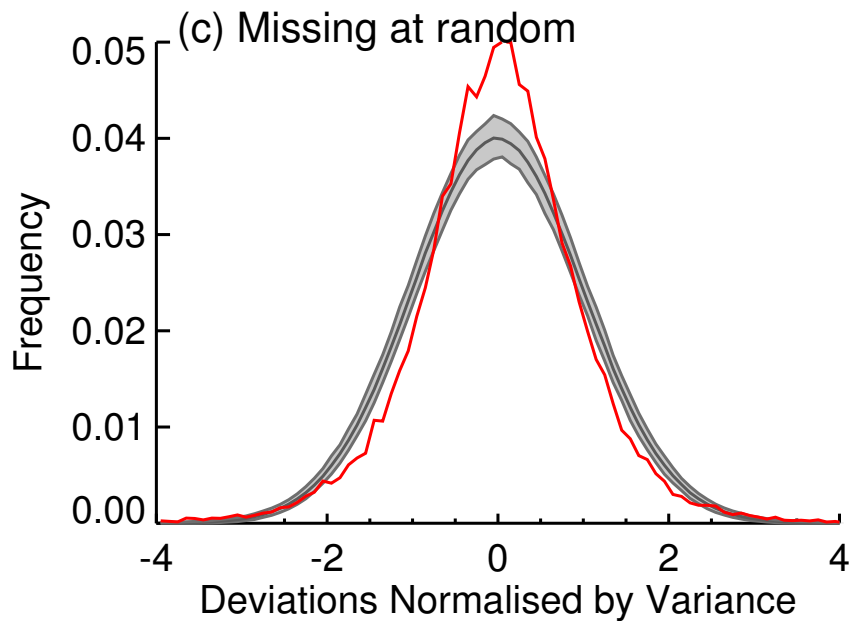
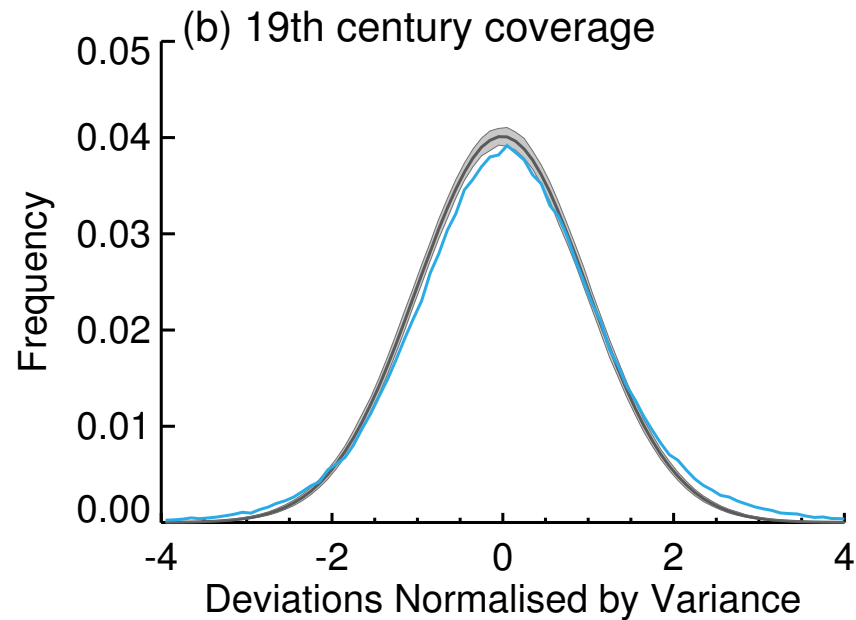
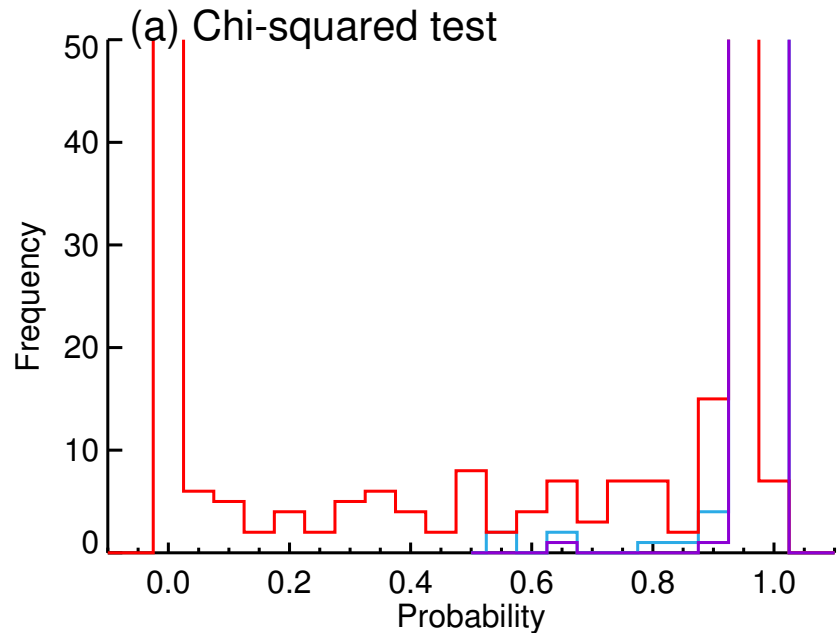


Figure 20.

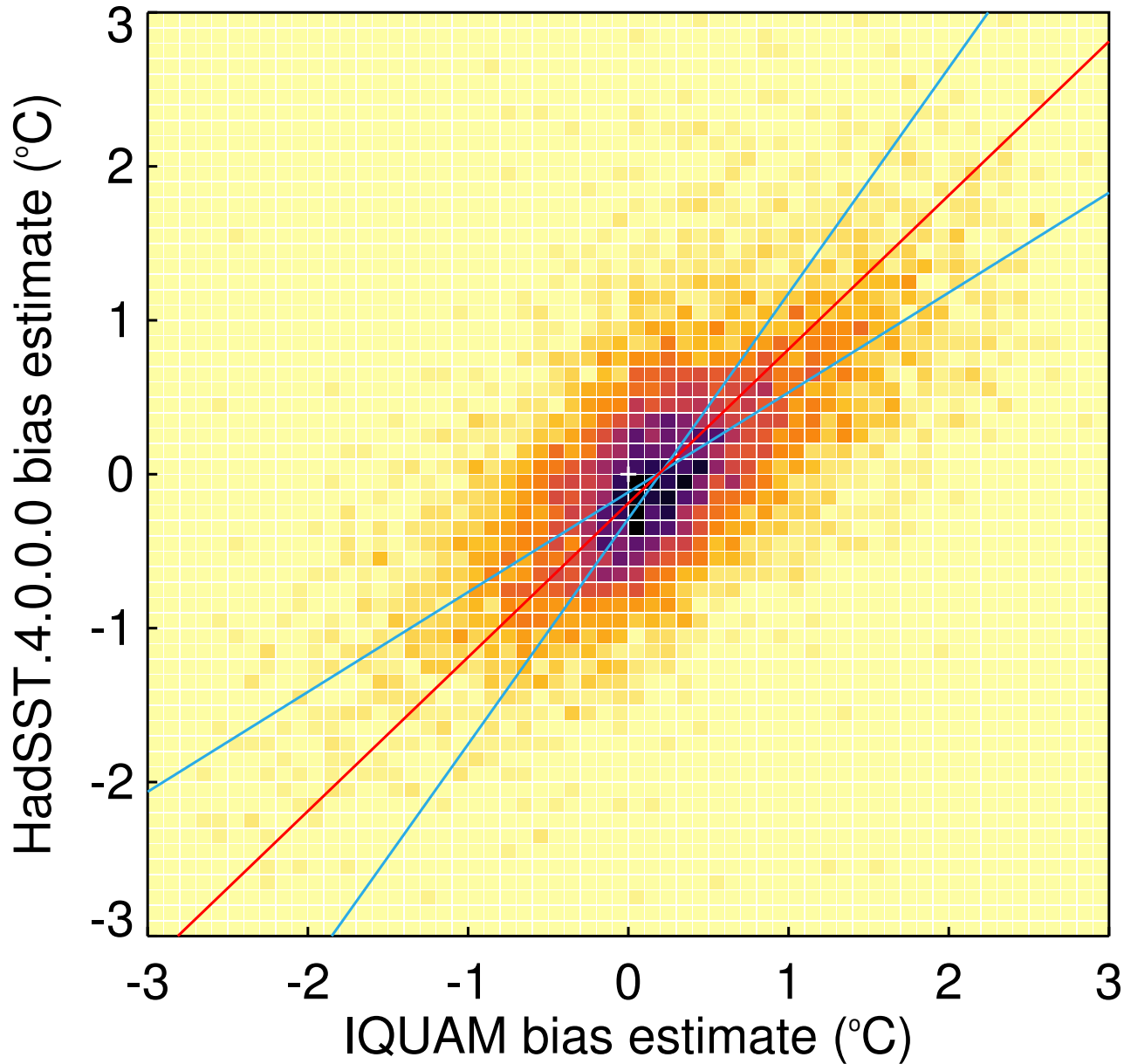


Figure 21.

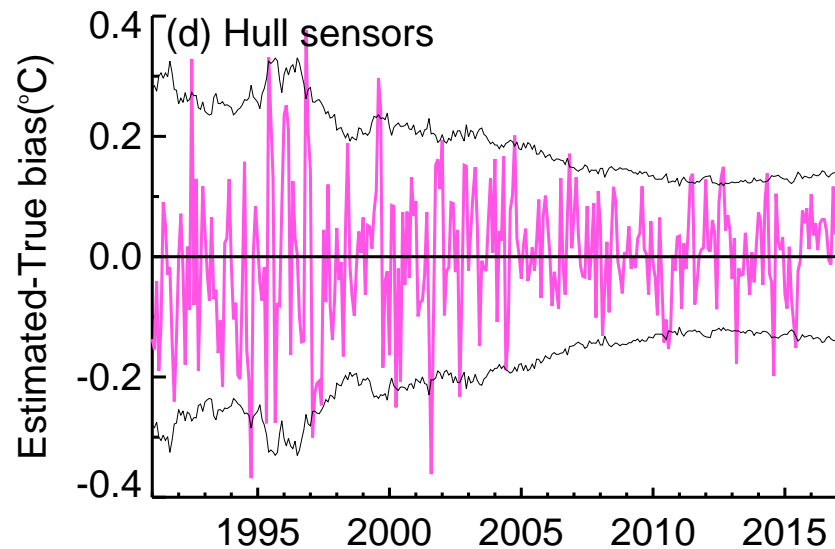
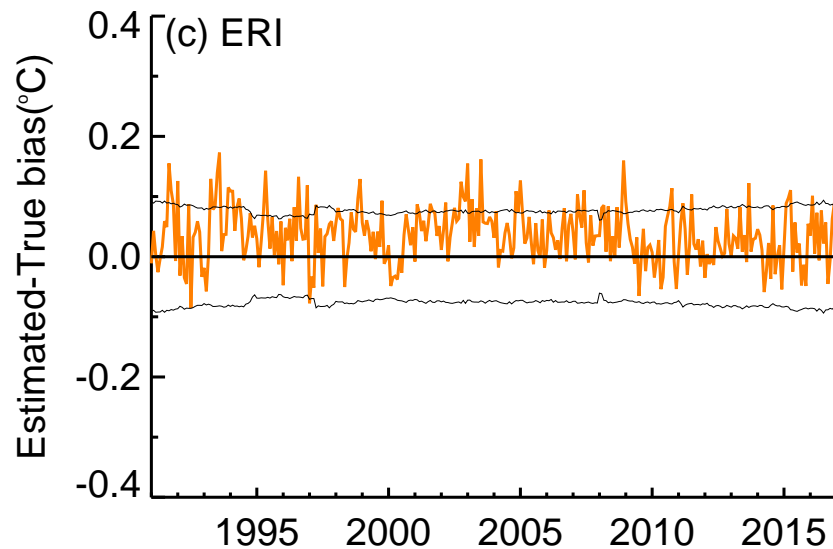
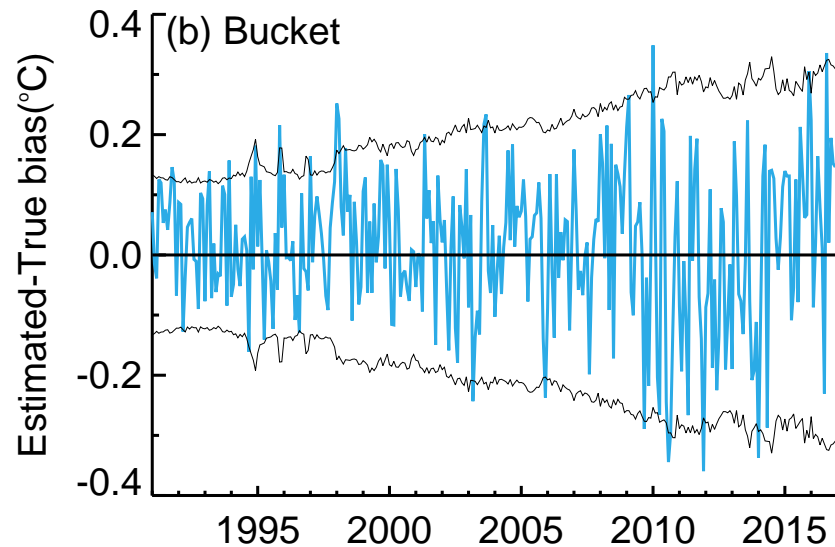
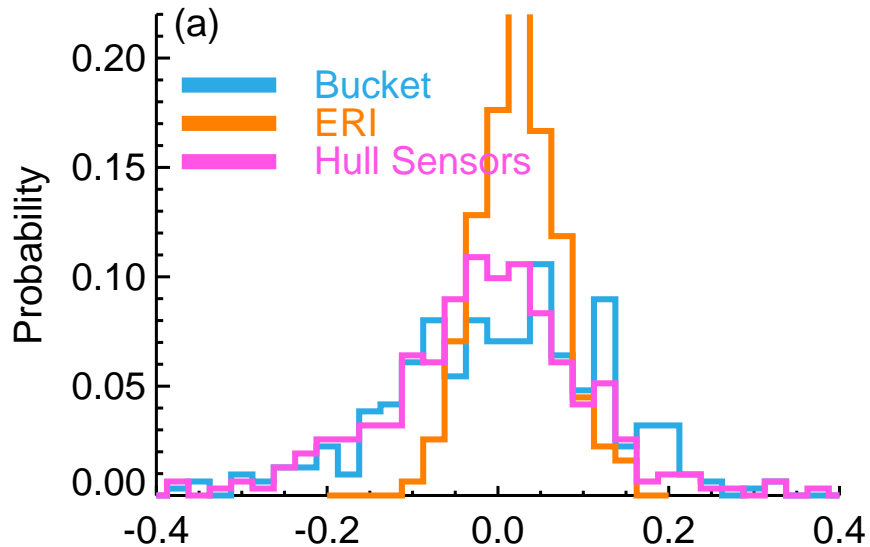


Figure 22.

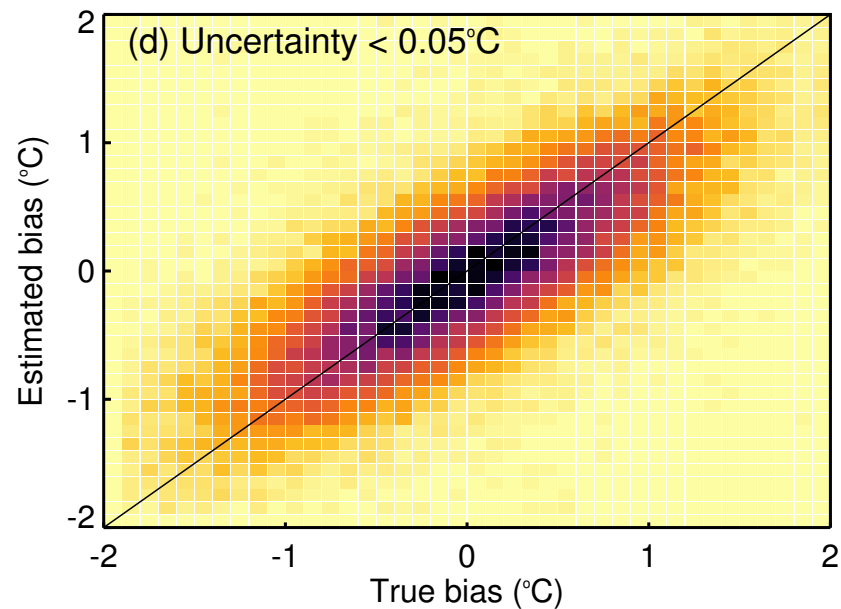
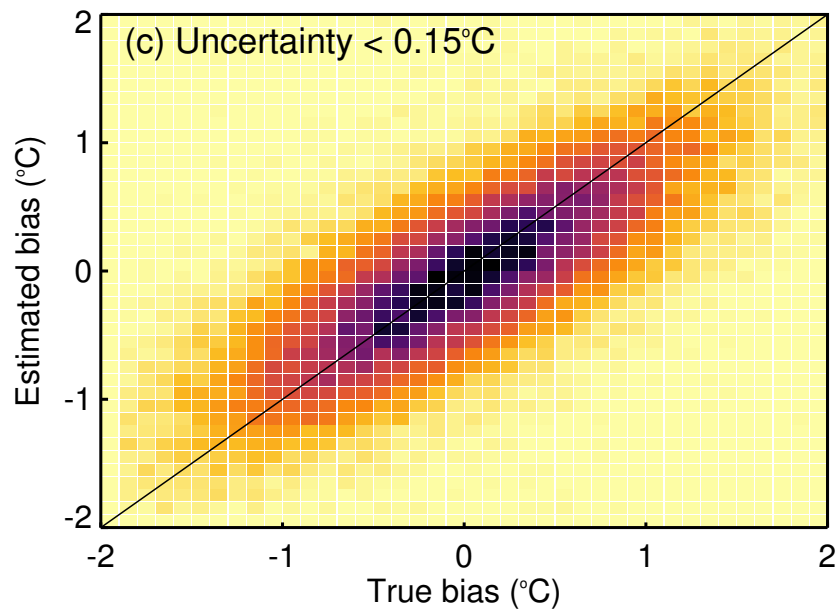
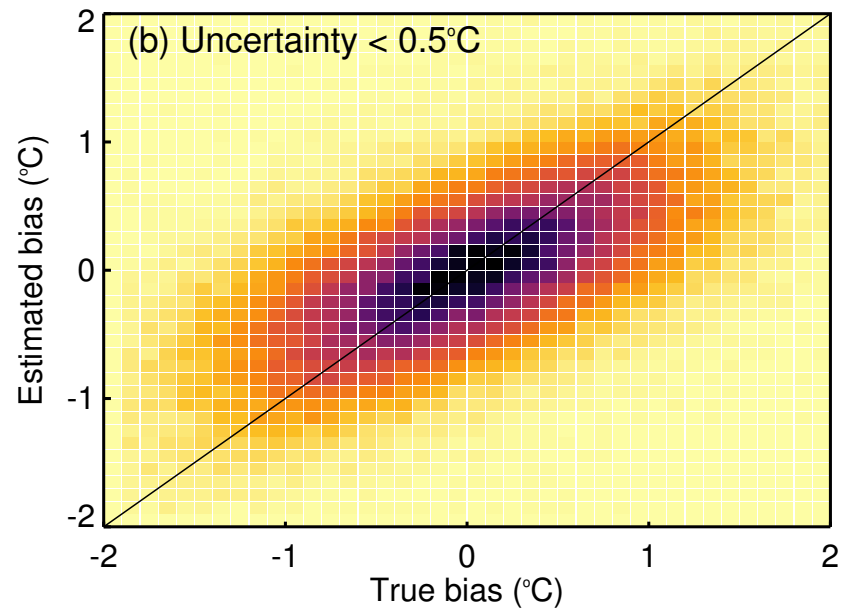
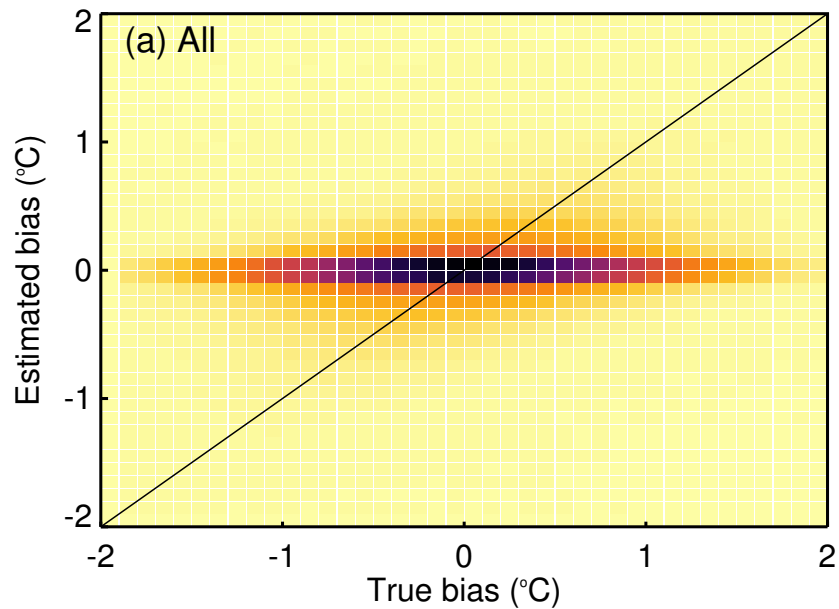


Figure 23.

